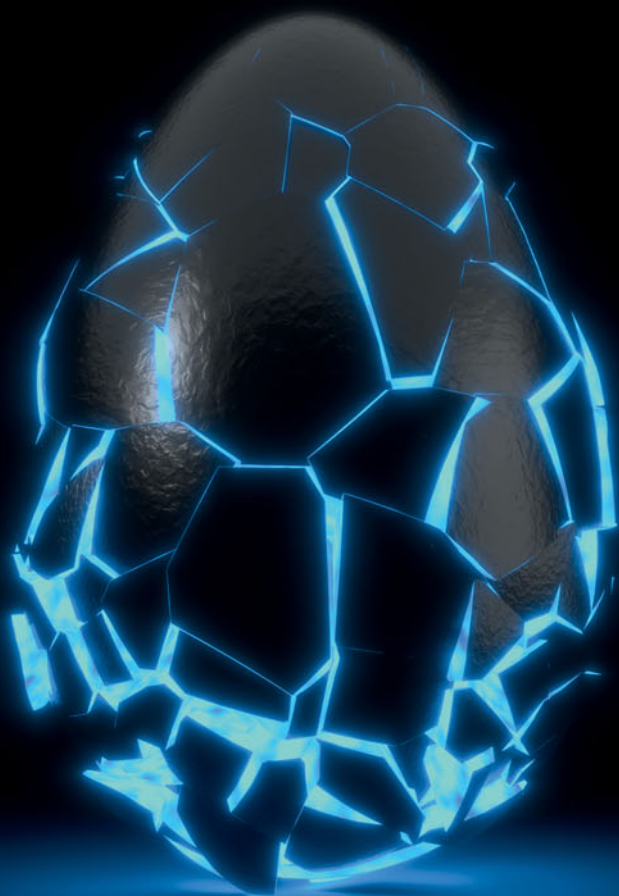


МАШИННОЕ ОБУЧЕНИЕ

ДЛЯ БИЗНЕСА И МАРКЕТИНГА



ИЛЬЯ КАЦОВ



Ilya Katsov

Introduction to Algorithmic Marketing:

Artificial Intelligence
for Marketing Operations

ИЛЬЯ КАЦОВ

МАШИННОЕ ОБУЧЕНИЕ

**ДЛЯ БИЗНЕСА
И МАРКЕТИНГА**



**Санкт-Петербург • Москва • Екатеринбург • Воронеж
Нижний Новгород • Ростов-на-Дону
Самара • Минск**

2019

ББК 32.813
УДК 004.8
К30

Кацов Илья

К30 Машинное обучение для бизнеса и маркетинга. — СПб.: Питер, 2019. — 512 с.: ил. — (Серия «IT для бизнеса»).

ISBN 978-5-4461-0926-5

Наука о данных становится неотъемлемой частью любой маркетинговой деятельности, и эта книга является живым портретом цифровых преобразований в маркетинге. Анализ данных и интеллектуальные алгоритмы позволяют автоматизировать трудоемкие маркетинговые задачи. Процесс принятия решений становится не только более совершенным, но и более быстрым, что имеет большое значение в постоянно ускоряющейся конкурентной среде.

«Эта книга — живой портрет цифровых преобразований в маркетинге. Она показывает, как наука о данных становится неотъемлемой частью любой маркетинговой деятельности. Подробно описывается, как подходы на основе анализа данных и интеллектуальных алгоритмов способствуют глубокой автоматизации традиционно трудоемких маркетинговых задач. Процесс принятия решений становится не только более совершенным, но и более быстрым, что важно в нашей постоянно ускоряющейся конкурентной среде. Эту книгу обязательно должны прочитать и специалисты по обработке данных, и специалисты по маркетингу, а лучше, если они будут читать ее вместе» (Андрей Себрант, директор по стратегическому маркетингу, Яндекс).

16+ (В соответствии с Федеральным законом от 29 декабря 2010 г. № 436-ФЗ.)

ББК 32.813
УДК 004.8

Права на издание получены по соглашению с Ilya Katsov. Все права защищены. Никакая часть данной книги не может быть воспроизведена в какой бы то ни было форме без письменного разрешения владельцев авторских прав.

Информация, содержащаяся в данной книге, получена из источников, рассматриваемых издательством как надежные. Тем не менее, имея в виду возможные человеческие или технические ошибки, издательство не может гарантировать абсолютную точность и полноту приводимых сведений и не несет ответственности за возможные ошибки, связанные с использованием книги. Издательство не несет ответственности за доступность материалов, ссылки на которые вы можете найти в этой книге. На момент подготовки книги к изданию все ссылки на интернет-ресурсы были действующими.

ISBN 978-0692989043 англ.
ISBN 978-5-4461-0926-5

© Ilya Katsov
© Перевод на русский язык ООО Издательство «Питер», 2019
© Издание на русском языке, оформление ООО Издательство «Питер», 2019
© Серия «IT для бизнеса», 2019
© Киселев А. Н., перевод на русский язык, 2019

Оглавление

Благодарности	12
Глава 1. Введение	13
1.1. Предмет алгоритмического маркетинга.....	14
1.2. Определение алгоритмического маркетинга	16
1.3. Исторические предпосылки и контекст	17
1.3.1. Онлайн-реклама: услуги и биржи.....	17
1.3.2. Авиакомпании: управление доходами	20
1.3.3. Наука маркетинга	22
1.4. Программные услуги.....	24
1.5. Кому адресована эта книга?.....	28
1.6. Итоги.....	29
Глава 2. Обзор предиктивного моделирования.....	30
2.1. Описательная, предиктивная и предписывающая аналитика	30
2.2. Экономическая оптимизация.....	31
2.3. Машинное обучение	34
2.4. Обучение с учителем	36
2.4.1. Параметрические и непараметрические модели.....	37
2.4.2. Оценка методом максимального правдоподобия	39
2.4.3. Линейные модели	40
2.4.4. Нелинейные модели.....	48
2.5. Обучение представлениям	53
2.5.1. Метод главных компонент.....	54
2.5.2. Кластеризация	61
2.6. Более специализированные модели	64
2.6.1. Теория потребительского выбора	64
2.6.2. Анализ выживаемости.....	71
2.6.3. Теория аукционов	81
2.7. Итоги	86

Глава 3. Продвижение и реклама	88
3.1. Среда.....	89
3.2. Бизнес-цели.....	93
3.2.1. Производители и ретейлеры	93
3.2.2. Затраты	94
3.2.3. Выгоды	95
3.3. Конвейер таргетирования	99
3.4. Моделирование и оценка отклика	101
3.4.1. Платформа моделирования отклика	102
3.4.2. Оценка отклика.....	106
3.5. Строительные блоки: таргетирование и модели ценности клиента	107
3.5.1. Сбор данных	108
3.5.2. Многоуровневое моделирование.....	110
3.5.3. RFM-моделирование.....	112
3.5.4. Моделирование предрасположенности	112
3.5.5. Сегментирование и персонализированное моделирование.....	122
3.5.6. Таргетирование с использованием анализа выживаемости	124
3.5.7. Моделирование пожизненной ценности	128
3.6. Проектирование и проведение кампаний.....	136
3.6.1. Цикл взаимодействий с клиентом.....	136
3.6.2. Кампании по продвижению продуктов	138
3.6.3. Многоступенчатые рекламные кампании.....	146
3.6.4. Кампании по удержанию.....	149
3.6.5. Кампании пополнения.....	152
3.7. Распределение ресурсов	153
3.7.1. Распределение между каналами.....	154
3.7.2. Распределение по целям	159
3.8. Онлайн-реклама	160
3.8.1. Среда.....	160
3.8.2. Цели и оценка.....	163
3.8.3. Таргетирование для модели CPA-LT	166
3.8.4. Оценка для случая с несколькими каналами	171
3.9. Оценка эффективности.....	174
3.9.1. Рандомизированные эксперименты	174
3.9.2. Неэкспериментальное исследование	181
3.10. Архитектура систем таргетирования	190
3.10.1. Сервер таргетирования.....	190
3.10.2. Платформа управления данными	192
3.10.3. Аналитическая платформа	192
3.11. Итоги	193

Глава 4. Поиск	196
4.1. Среда.....	197
4.2. Бизнес-цели.....	200
4.2.1. Метрики релевантности	202
4.2.2. Средства управления продвижением	207
4.2.3. Метрики качества службы поиска	209
4.3. Строительные блоки: соответствие и ранжирование.....	210
4.3.1. Лексическое соответствие.....	211
4.3.2. Логический поиск и поиск по фразе.....	213
4.3.3. Нормализация и стемминг.....	214
4.3.4. Ранжирование и модель векторного пространства	216
4.3.5. Модель оценки $TF \times IDF$	219
4.3.6. Оценка с использованием n-грамм	223
4.4. Смешивание сигналов релевантности	224
4.4.1. Поиск по нескольким полям	225
4.4.2. Проектирование и регулировка сигналов	227
4.4.3. Проектирование конвейера смешивания сигналов	234
4.5. Семантический анализ	237
4.5.1. Синонимы и иерархии	238
4.5.2. Векторные представления слов	241
4.5.3. Латентно-семантический анализ	243
4.5.4. Вероятностное тематическое моделирование.....	251
4.5.5. Вероятностный латентно-семантический анализ	252
4.5.6. Латентное размещение Дирихле	257
4.5.7. Модель Word2Vec	259
4.6. Методы поиска для продвижения.....	267
4.6.1. Комбинаторный фразовый поиск.....	269
4.6.2. Контролируемое снижение точности	273
4.6.3. Вложенные сущности и динамическая группировка	274
4.7. Настройка релевантности.....	278
4.7.1. Обучение ранжированию	279
4.7.2. Обучение ранжированию на неявной обратной связи	285
4.8. Архитектура служб поиска товаров.....	289
4.9. Итоги	291
Глава 5. Рекомендации	294
5.1. Среда.....	296
5.1.1. Свойства рейтингов клиентов	298
5.2. Бизнес-цели.....	300
5.3. Оценка качества	302

5.3.1. Точность прогнозирования.....	303
5.3.2. Точность ранжирования.....	305
5.3.3. Новизна	307
5.3.4. Серендипность.....	307
5.3.5. Разнообразие.....	308
5.3.6. Охват	308
5.3.7. Роль экспериментов	309
5.4. Обзор методов рекомендаций.....	310
5.5. Фильтрация по содержанию	312
5.5.1. Метод ближайших соседей.....	316
5.5.2. Наивный байесовский классификатор	317
5.5.3. Проектирование признаков для фильтрации по содержанию	323
5.6. Введение в совместную фильтрацию	325
5.6.1. Базовые оценки	328
5.7. Совместная фильтрация на основе близости.....	331
5.7.1. Совместная фильтрация по близости пользователей.....	333
5.7.2. Совместная фильтрация по близости элементов.....	339
5.7.3. Сравнение методов на основе близости пользователей и элементов	341
5.7.4. Методы на основе близости как задача регрессии	342
5.8. Совместная фильтрация на основе моделей	348
5.8.1. Адаптация регрессионных моделей для предсказания рейтингов.....	349
5.8.2. Наивная байесовская совместная фильтрация	351
5.8.3. Модели скрытых факторов	356
5.9. Гибридные методы.....	371
5.9.1. Переключение	372
5.9.2. Смешивание.....	373
5.9.3. Расширение признаков	379
5.9.4. Варианты представления гибридных рекомендаций	382
5.10. Контекстные рекомендации	383
5.10.1. Многомерная основа	384
5.10.2. Контекстно-зависимые методы рекомендаций	386
5.10.3. Модели рекомендаций с учетом времени	389
5.11. Неперсонализированные рекомендации.....	394
5.11.1. Типы неперсонализированных рекомендаций	394
5.11.2. Рекомендации с использованием ассоциативных правил.....	396
5.12. Многоцелевая оптимизация	400
5.13. Архитектура систем рекомендаций.....	404
5.14. Итоги	406

Глава 6. Ценообразование и ассортимент	409
6.1. Среда.....	410
6.2. Влияние ценообразования	413
6.3. Цена и стоимость.....	414
6.3.1. Ценовые границы.....	415
6.3.2. Субъективная ценность.....	417
6.4. Цена и спрос.....	419
6.4.1. Линейная кривая спроса	421
6.4.2. Кривая спроса с постоянной эластичностью.....	422
6.4.3. Логит-кривая спроса	423
6.5. Базовые структуры цен.....	425
6.5.1. Цена за единицу	426
6.5.2. Сегментация рынка.....	428
6.5.3. Комплексное ценообразование	433
6.5.4. Пакетирование.....	437
6.6. Прогнозирование спроса.....	441
6.6.1. Модель спроса для оптимизации ассортимента	443
6.6.2. Модель спроса для сезонных продаж	446
6.6.3. Прогнозирование спроса с учетом исчерпания запасов.....	449
6.7. Оптимизация цен	452
6.7.1. Ценовая дифференциация	453
6.7.2. Динамическое ценообразование.....	462
6.7.3. Персонализированные скидки	472
6.8. Распределение ресурсов	475
6.8.1. Среда.....	475
6.8.2. Распределение с двумя классами	479
6.8.3. Распределение с несколькими классами.....	482
6.8.4. Эвристические решения для нескольких классов	484
6.9. Оптимизация ассортимента.....	486
6.9.1. Оптимизация планировки магазина.....	487
6.9.2. Управление категориями.....	490
6.10. Архитектура систем управления ценами	496
6.11. Итоги	498
Приложение. Распределение Дирихле	500
Библиография.....	504

Сейчас, когда первенство отдается потребителю и бренды с ретейлерами хватаются за малейшие проявления внимания, идет жесткая конкуренция за данные и возможность их использования для целеполагания, привлечения и удержания клиентов. Книга является руководством, рассказывающим, как выстоять в этой борьбе. Она будет полезна специалистам по маркетингу и разработчикам, проведет их по всей маркетинговой цепочке и покажет, как оцифровать ее. Исчерпывающий и незаменимый справочник для всех, кто встал на путь алгоритмического маркетинга.

Али Бухуч, технический директор Sephora Americas

Сейчас возможно все. Эта книга подводит практический фундамент под понятия, которые всего несколько лет назад отвергались как чистая теория. В ней дается четкое обоснование того, что лучшие маркетологи понимают на интуитивном уровне, но не могут выразить словами. В элегантной математической форме формулируются важные отношения, неуловимые для традиционного бизнес-моделирования. Книге не нужны оправдания из-за отсутствия развернутых примеров — большую часть мира нельзя представить линейно, пользуясь лишь несколькими измерениями и не допуская неопределенности. Вместо этого книга устанавливает строгие рамки, помогающие лучше понять реальные явления. Она написана не для исследователей данных и не для маркетологов, а скорее для тех и других вместе! Только партнерство таких специалистов позволит получить настоящую отдачу. С этой книги должно начаться такое партнерство.

Эрик Колсон, руководитель методического отдела Stitch Fix

Эта книга — живой портрет цифровых преобразований в маркетинге. Она показывает, как наука о данных становится неотъемлемой частью любой маркетинговой деятельности. В книге подробно описывается, как подходы на основе анализа данных и интеллектуальных алгоритмов способствуют глубокой автоматизации традиционно трудоемких маркетинговых задач. Процесс принятия решений становится не только более совершенным, но и более быстрым, что имеет большое значение в нашей постоянно ускоряющейся конкурентной среде. Эту книгу обязательно должны прочитать и специалисты по обработке данных, и специалисты по маркетингу, а лучше, если они будут читать ее вместе.

Андрей Себрант, директор по стратегическому маркетингу, «Яндекс»

Эта книга содержит комплексный план, как целиком оцифровать маркетинговую деятельность вашей компании. Начиная с концептуального описания архитектуры будущего цифрового маркетинга, она углубляется в детальный анализ передовых подходов в каждой отдельной области маркетинговых операций. Книга даст руководителям, менеджерам среднего звена и специалистам по обработке данных в вашей организации набор конкретных, практичных и поэтапных рекомендаций, как, приступая с сегодняшнего дня, начать генерировать все лучшие идеи и решения.

Виктория Лившиц, основатель и технический директор Grid Dynamics

Книга предлагает практикующим маркетологам ценные рецепты использования продвинутых методов машинного обучения и науки о данных, помогающих понять поведение клиентов, персонализировать предложения, оптимизировать систему поощрений и управлять вовлеченностью, способствуя тем самым созданию нового поколения платформ анализа данных для маркетинговых систем.

Кира Макагон, директор по инновациям, RingCentral; серийный предприниматель, основатель RedAriL и Octane

Сегодня почти все руководители, отвечающие за коммерческую часть, понимают концептуальную важность анализа данных и машинного обучения, тем не менее задача внедрения актуальных конкурентных решений, основанных на науке о данных, остается довольно сложной. К числу трудностей, с которыми сталкиваются специалисты, занимающиеся цифровым маркетингом, можно причислить нехватку талантливых исследователей данных и сложность адаптации к конкретным условиям академических моделей, универсального открытого программного обеспечения и алгоритмов. Книга Ильи Кацова основана на богатом опыте, накопленном в Grid Dynamics, в области разработки инновационных, но практичных решений цифрового маркетинга для крупных организаций, которые помогают им успешно конкурировать, идти в ногу со временем и адаптироваться к новой эпохе в анализе данных.

Эрик Бенаму, основатель и генеральный партнер Benhamou Global Ventures; бывший генеральный директор и президент 3Com и Palm

Благодарности

Эта книга была бы невозможна без поддержки и помощи многих людей. Я очень благодарен моим коллегам и друзьям: Али Бухучу (Ali Bouhouch), Максиму Мартынову (Max Martynov), Дэвиду Нейлору (David Naylor), Пенелопе Конлон (Penelope Conlon), Сергею Трюберу (Sergey Tryuber), Денису Копыченко (Denys Kopyuchenko) и Вадиму Козырькову (Vadim Kozyrkov), которые прочитали рукопись и поделились своими мнениями. Отдельное спасибо Константину Перикову (Konstantin Perikov), представившему массу содержательных предложений по поисковым службам, а также оказавшему помощь с некоторыми примерами.

Я очень обязан Игорю Яговому (Igor Yagovoy), Виктории Лившиц (Victoria Livschitz), Леонарду Лившицу (Leonard Livschitz) и Эзре Бергеру (Ezra Berger) за поддержку этого проекта и помощь с изданием. Наконец, я хочу сказать большое спасибо Кэтрин Райт (Kathryn Wright) — моему редактору, которая помогла мне превратить рукопись в законченный продукт.

1

Введение

В 1888 году Винсент Ван Гог, в ту пору малоизвестный голландский художник, написал своему брату Тео, что *художником будущего может стать лишь невиданный еще колорист*. Не вдаваясь в аспекты изобразительного искусства, поражает и восхищает сам способ, каким Ван Гог ставит вопрос о художниках будущего и дает на него ответ. Ван Гог, несомненно, был прав, предвидя, что художники грядущего века будут развивать навыки, прежде неведомые, и изменят путь развития искусства. А что, если тот же вопрос задать в отношении практикующих маркетологов, живущих в эпоху цифровых средств массовой информации и изобилия данных? Кем будет маркетолог будущего? Будет ли он художником связей с клиентами? Статистиком, невиданным прежде? Программистом, создающим маркетинговые системы? Экспертом в экономическом моделировании?

Историю маркетинга можно рассматривать как эволюцию принципов, приемов и методов оптимизации бизнеса. Всегда считалось, что к этой проблеме оптимизации можно подойти с научной точки зрения и применить строгие математические методы к широкому кругу маркетинговых задач. Однако приверженцы таких методов неизбежно сталкивались с проблемами, связанными с неполнотой данных, сложностями маркетинга в реальной жизни, негибкостью бизнес-процессов и ограничениями ПО. Особенно острые проблемы наблюдаются в областях, где требуется принимать далеко идущие стратегические решения, где человеческое суждение зачастую является единственным жизнеспособным решением для практического применения.

Развитие каналов цифрового маркетинга изменило игру и создало среду, требующую принятия миллионов *микрорешений*, что просто невозможно без интеллектуального программного обеспечения и алгоритмов. Целевые рекламные акции, динамическое ценообразование в обычных и интернет-магазинах, службы поиска и подбора рекомендаций электронной коммерции, онлайн-реклама — все это требует применения продвинутых методов экономического моделирования, науки

о данных и разработки программного обеспечения для реализации потенциала цифровой среды. Например, этот потенциал нельзя реализовать полностью без учета индивидуальных потребностей миллионов клиентов, что в свою очередь требует принятия миллионов уникальных решений. Кроме того, вездесущая цифровая среда и мобильные устройства позволяют клиентам пройти маркетинговую воронку от поиска до покупки в считанные секунды, в любом месте и в любое время, и эта распространенная *ситуация микромомента* также требует принятия маркетинговых решений за микросекунды. Такой характер среды порождает проблему построения маркетинговых систем, которые принимают решения и действуют на беспрецедентном уровне автономности, производя широкий и глубокий анализ. На основе анализа данных в некоторых случаях возможно не только принятие отдельных решений, но также планирование, выполнение и оптимизация целых бизнес-процессов с привлечением автоматизированных программных систем.

Автоматизации могут подвергаться разные аспекты маркетинга, включая экономику, управление, статистику и анализ, однако создатели таких систем должны расставить все эти части как единый набор методов, которые можно эффективно реализовать в программном обеспечении для достижения бизнес-целей. Руководить современным проектом в области маркетинговых технологий — все равно что дирижировать оркестром, включающим разнообразные инструменты, и заставлять их звучать в унисон. Именно с этой точки зрения мы будем рассматривать маркетинг на протяжении всей книги, обобщать богатый опыт, накопленный за последние десятилетия разработчиками маркетинговых систем в ретейле, онлайн-рекламе и других отраслях, и исследовать руководящие теоретические принципы. Следует отметить, что мы сознательно ориентируемся не на академические и теоретические исследования, а на результаты, представленные практикующими специалистами, доказавшими эффективность в бизнес-решениях. К счастью, число методов, моделей и архитектур, опубликованных такими практиками, достаточно велико и иногда они описываются с большим количеством подробностей. Одна часть этих публикаций сосредоточена в основном на технологиях и аспектах реализации, другая больше внимания уделяет математическому моделированию, оптимизации и экономико-математическим методам. В действительности оба аспекта важны для создания и функционирования успешной маркетинговой системы. Многие из опубликованных результатов также основаны на моделях, разработанных академическими исследователями в области научного маркетинга.

1.1. Предмет алгоритмического маркетинга

Одно из традиционных определений маркетинга описывает его как деятельность по определению продуктов и услуг, предлагаемых компанией, и их представле-

нию настоящим или потенциальным клиентам. Эту деятельность можно разбить на несколько потоков, которые обычно описываются как вариации следующих категорий [McCarthy, 1960]:

- Продукт — анализ маркетинговых возможностей, планирование линеек продуктов и их характеристик, планирование ассортимента.
- Продвижение — все методы коммуникации между компанией и ее клиентами: реклама, рекомендации, обслуживание клиентов и др.
- Цена — стратегии ценообразования, включая объявленные цены, скидки и изменение цен с течением времени.
- Распространение — исторически под этим понимается процесс предоставления продукта или услуги конечному пользователю через разные каналы. Более поздние интерпретации подчеркивают роль открытия продукта и удобство его покупки, при этом отмечается снижение важности роли распространения с увеличением числа каналов цифрового маркетинга [Lauterborn, 1990].

Эта классификация широко известна как *маркетинг-микс*, или *модель 4P¹*. Этот микс можно рассматривать как набор переменных, которые могут контролировать-ся маркетологом или маркетинговым ПО для влияния на положение продуктов и брендов на рынке. Каждый компонент маркетинг-микса представляет широкую область, которую можно рассматривать и изучать с разных сторон. Предмет алгоритмического маркетинга проще понять, выделив два аспекта маркетинговой деятельности: *стратегию* и *процесс*. Под термином *стратегия* в данном случае понимаются долгосрочные бизнес-решения высокого уровня, которые определяют конкурентные преимущества компании и общее направление ее маркетинговых процессов. Например, ретейлер должен определить свой целевой рынок, порядок обслуживания клиентов и линейки продуктов как части бизнес-стратегии. *Процесс* — это реализация стратегии, ориентированной на тактические решения, обеспечивающие непрерывное функционирование компании. Продолжая пример с ретейлером, высокоуровневые стратегии ценообразования и продвижения требуют многочисленных решений, определяющих, как выбирать потребителей кампаний по продвижению или как с течением времени должны меняться цены на отдельные продукты.

Несмотря на то что сферы стратегических и тактических процессов не имеют строгого деления и между ними нет четкой границы, можно утверждать, что стратегия больше сосредоточена на исследовании, анализе и планировании с участием

¹ *Product* — товар или услуга; *Price* — цена, наценки, скидки; *Promotion* — продвижение, реклама, пиар, стимулирование сбыта; *Place* — месторасположение торговой точки, каналы распределения, персонал продавца. — *Примеч. ред.*

человеческого суждения, тогда как процесс больше сосредоточен на выполнении, принятии микрорешений и, самое главное, на автоматизации. Это делает процесс особенно привлекательным для нашего исследования, хотя и стратегию, и процесс можно описать с точки зрения науки о данных и извлечь выгоду из методов на основе анализа данных. Подводя краткий итог, можно сказать, что предметом алгоритмического маркетинга в основном являются процессы в четырех областях маркетинг-микса, а также автоматизация этих процессов с использованием методов на основе анализа данных и эконометрики.

1.2. Определение алгоритмического маркетинга

Мы определяем *алгоритмический маркетинг* как процесс маркетинга, автоматизированный до такой степени, что им можно управлять, задавая бизнес-цель в маркетинговой программной системе. Это означает, что маркетинговая система должна быть достаточно интеллектуальной, чтобы понимать высокоуровневые цели, такие как привлечение новых клиентов или максимизация доходов, чтобы спланировать и выполнить последовательность бизнес-мероприятий, таких как рекламная кампания или корректировка цен, для достижения цели и обучаться на результатах, чтобы корректировать и оптимизировать мероприятия при необходимости. Этот принцип проиллюстрирован на рис. 1.1. В этой книге мы также будем использовать термин *программный* для обозначения высокоавтоматизированных маркетинговых программных систем и служб, и в большинстве контекстов использовать термины *алгоритмический* и *программный* взаимозаменяемо.

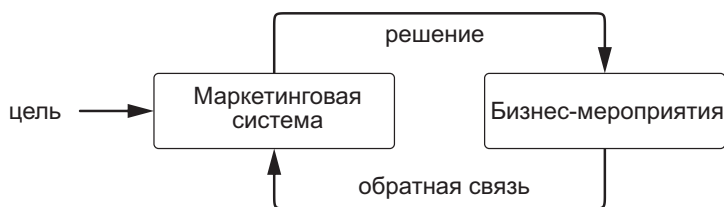


Рис. 1.1. Концептуальное представление экосистемы алгоритмического маркетинга

Было бы идеально, если бы программная система была полностью автоматизированной и автономной, но мы не считаем это принципиальной целью или обязательным требованием. Напротив, программная система обычно поддерживается множеством специалистов, включая исследователей данных, инженеров и аналитиков, которые разрабатывают и корректируют модели и алгоритмы с целью увеличения эффек-

тивности системы и расширения ее возможностей. Также она может использовать результаты стратегического анализа и планирования, полученные с применением непрограммных методов, и, возможно, в связи с некоторыми другими задачами. Однако способность системы понимать бизнес-цели и получать измеримые результаты очень важна. При этом нужно помнить об ограничениях и опасностях автоматизации в маркетинге. Во многих прикладных применениях уместнее рассматривать программные системы как интеллектуальные инструменты, помогающие маркетологам эффективнее достигать желаемого, а не как их замену.

1.3. Исторические предпосылки и контекст

Между алгоритмическим и неалгоритмическим маркетингом нет четкой границы. Более того, проводить такую границу в некотором смысле недопустимо, потому что алгоритмические системы — это лишь ответ на старые вопросы, а не постановка новых маркетинговых вопросов. Однако очевидно, что уровень автоматизации маркетинга сильно разнится по отраслям. Это свидетельствует о том, что в одних условия более благоприятны, чем в других. И наоборот, успешное внедрение передовых алгоритмических методов может существенно преобразить отрасль и создать еще более благоприятные условия для дальнейшего развития. Анализ такой благоприятной среды является естественной отправной точкой для понимания алгоритмического маркетинга. Давайте вкратце рассмотрим несколько бизнес-сценариев, заложивших основы алгоритмического маркетинга, и определим их общие черты и закономерности, позволившие применить системный подход.

1.3.1. Онлайн-реклама: услуги и биржи

История интернет-рекламы началась 3 мая 1978 года, когда первое рекламное электронное письмо было отправлено 400 пользователям компьютерной сети ARPANET, развернутой в то время в четырех местах: в университете штата Юта, Калифорнийском университете в Лос-Анжелесе, Калифорнийском университете в Санта-Барбаре и в исследовательском центре Стэнфордского университета. Через пятнадцать лет, в 1993 году, когда ARPANET переросла в интернет с множеством мультимедийных веб-сайтов, возник рынок рекламных баннеров. Первоначально этот новый рынок полагался на прямую продажу мест для баннеров держателями веб-сайтов рекламодателям, но этот подход быстро утратил свою эффективность со всплеском роста веб-сайтов. Рекламодателям стало сложно, а иногда невозможно проводить рекламные кампании и управлять бюджетами при наличии тысяч издателей. С другой стороны, издателям нужен был надежный и централизованный способ продажи рекламных площадей.

Вызов был принят рекламными сетями, выступившими в роли брокеров между издателями и рекламодателями. Проект DoubleClick, основанный в 1996 году, предложил платформу, которая позволила рекламодателям проводить рекламные кампании в обширной сети веб-сайтов, динамически вносить коррективы в соответствии с их эффективностью и оценивать отдачу инвестиций. Это создало идеальную среду для автоматического принятия решений, благодаря возможности динамически оценивать результат и вносить коррективы. Однако это был не настоящий программный подход.

В то же время системы онлайн-поиска из всех сил стремились усовершенствовать свои возможности по доставке рекламы. Рекламодатели платили за количество показов их рекламы в поисковых системах — согласно модели «цена за тысячу показов (Cost Per Thousand, или Cost per Mille, CPM) — по аналогии с моделью показа баннеров. Этот подход был негибким с точки зрения ценообразования, приводивший к потерям доходов поисковых систем, и малоэффективным с точки зрения целевой аудитории, потому что показ нерелевантных объявлений никак не отражался на стоимости. Прорыв случился в 1998 году, когда поисковая система GoTo.com представила автоматизированную модель аукциона с двумя инновационными особенностями:

- Рекламодатели делали ставки за появление в первых строках результатов, возвращаемых в ответ на поисковые запросы.
- Рекламодатели платили за количество переходов по рекламным ссылкам, а не за количество показов.

Модель платы за клик (Per-Pay-Click, PPC) способствовала повышению и доходности, и релевантности, так как рекламодатели, готовые платить за верхние строчки в результатах для конкретных запросов, обычно предлагали более релевантные ресурсы. Эта модель была взята на вооружение компанией Google в 2002 году с одним важным усовершенствованием: рекламные объявления выбирались, исходя из оценки доходности, а не величины ставки. Поисковая система Google измеряла удельную стоимость щелчка для каждой рекламной ссылки как соотношение щелчков и показов, а ожидаемая доходность оценивалась как

$$\text{доходность} = \text{величина ставки} \times \text{частота переходов по ссылке}.$$

Это уже был программный метод самообучения, оптимизировавший бизнес-цель с точки зрения доходности и релевантности, потому что величина частоты переходов по ссылке обычно невелика для нерелевантной рекламы, поэтому даже рекламодатели с большим бюджетом не могли заполнить рекламные площади.

Траектории рекламных сетей и поисковых систем сошлись в 2007–2009 годах с внедрением аукционной модели по всем направлениям. Рекламодателей и издателей

связали рекламные биржи, которые принимали ставки в режиме реального времени для показа отдельных объявлений. Началась новая эра аукционов рекламных объявлений в реальном времени (Real-Time Bidding, RTB). Появление RTB-бирж дало толчок развитию программных инструментов для рекламодателей — платформ управления данными (Data-Management Platform, DMP) и платформ для стороны рекламной сети (Demand-Side Platform, DSP), — которые обеспечили возможность сбора данных о поведении пользователей интернета и делать ставки на RTB-биржах, в зависимости от склонности данного пользователя откликаться на рекламу. Успех RTB оказался впечатляющим: доходы от продажи рекламы компанией DoubleClick (приобретенной компанией Google к тому времени) через RTB выросли с 8 % в январе 2010 до 68 % в мае 2011 [Google Inc., 2011].

Размышляя над историей RTB, можно сделать вывод, что одним из самых заметных достижений перевода рекламы на программные рельсы является создание инфраструктуры, позволившей владельцам баз данных с информацией о потребителях — издателям веб-контента — предоставить маркетинговые услуги сторонам, ограниченным в своих возможностях взаимодействовать с потребителями, — рекламодателям, рекламирующим свои продукты и услуги. Такая инфраструктура, находящаяся между издателями и рекламодателями, как правило, поддерживается независимой стороной и включает:

- *Рекламные службы*, позволяющие рекламодателям проводить рекламные кампании с использованием ресурсов издателя. Обычно эти услуги используются для установления связи между несколькими рекламодателями и издателями и напоминают торговые площадки, где продаются и покупаются ресурсы, часто на основе ставок.
- *Службы данных*, собирающие и хранящие информацию о потребителях, которая поступает от издателей, рекламодателей и третьих сторон. Рекламные службы используют эти данные для проведения рекламных кампаний и автоматического принятия решений о доставке рекламы в режиме реального времени.

Позже эта тенденция стала распространяться на другие отрасли. Другие владельцы баз данных с информацией о потребителях, например ретейлеры и операторы мобильной связи, тоже искали эффективный способ коммерциализации своих данных и связей с потребителями, а другие поставщики услуг, такие как банки, производители товаров и страховые компании, хотели больше информации о своих потребителях и иметь больше каналов связи с ними. Например, производитель товаров массового потребления может использовать каналы ретейлера, такие как магазины и веб-сайты электронной коммерции, чтобы предложить персональные скидки потребителям для продвижения новых товаров и увеличения их доли на рынке.

В результате рекламные службы и службы данных начали трансформироваться в более общую модель, изображенную на рис. 1.2, которая представляет собой многоцелевой рынок услуг и данных, объединяющий участников из разных отраслей. Спектр услуг, предлагаемых на таком рынке, может выходить далеко за рамки рекламы, охватывая такие области, как кредитование и страхование. Неравномерность этой среды, где постоянно обрабатываются сторонние данные, часто в режиме реального времени, обуславливает высокую сложность потоков данных и оперативных решений, из-за чего решение проблемы возможно, пожалуй, только с привлечением программных методов.

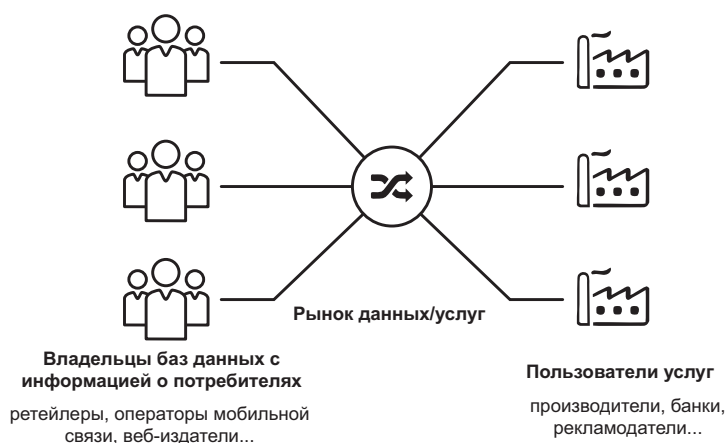


Рис. 1.2. Модель рынка данных и услуг

1.3.2. Авиакомпании: управление доходами

Онлайн-реклама и рынок данных — пожалуй, самые известные и успешные примеры в истории программного маркетинга, но точно не единственные. Онлайн-реклама в какой-то степени оказалась совершенно новой средой, предлагающей беспрецедентные возможности для решения задач, не имеющих аналогов. Однако программные методы с успехом могут применяться в более традиционных областях. Рассмотрим ситуацию, которая по совпадению стала развиваться в том же 1978 году, когда было отправлено первое рекламное электронное письмо.

С 1938 года Федеральный совет по гражданской авиации (Civil Aeronautics Board) регулировал работу всего федерального воздушного транспорта США, определяя расписание движения, маршруты и тарифы на основе стандартных цен и целевых показателей рентабельности авиакомпаний. Закон об отмене регулирования авиаперевозок от 1978 года устранил многие виды контроля и дал авиакомпаниям

свободно изменять цены и маршруты. Это открыло дверь перед бюджетными компаниями, которые первыми разработали более простые модели предоставления услуг без излишеств и уменьшили трудозатраты. Одним из самых ярких примеров стала компания People Express, основанная в 1981 году и предложившая тарифы на 70 % ниже, чем крупные авиакомпании.

Бюджетные перевозчики привлекли новые категории путешественников, которые до этого редко пользовались воздушным транспортом: студентов, приезжающих домой на каникулы, отдыхающих, выезжающих на несколько дней, и многих других. В 1984 году People Express отчиталась о доходе в 1 миллиард долларов США, а ее чистая прибыль составила 60 миллионов [Talluri and Van Ryzin, 2004]. Появление бюджетных авиакомпаний превратилось в угрозу крупным перевозчикам, у которых практически не было шансов выиграть ценовую войну. Кроме того, старые авиакомпании не могли позволить себе потерять высокодоходных бизнес-путешественников в погоне за мало доходным рынком.

Решение было найдено в American Airlines. Во-первых, они признали, что непроданные места можно использовать для конкуренции по цене с бюджетными перевозчиками, потому что расчетная стоимость таких мест все равно стремится к нулю. Но появилась другая проблема — как предотвратить покупку билетов по сниженным ценам бизнес-путешественниками? Для этого было решено ограничить льготное предложение; например, билеты со скидкой должны были приобретаться не менее чем за три недели до вылета и не подлежали возврату. Проблема заключалась в том, что избыток свободных мест существенно менялся между рейсами, и оптимального их распределения можно было достичь только с использованием динамической оптимизации. В 1985-м, после нескольких лет разработки, в American Airlines была запущена система оптимизации динамического распределения ресурсов и услуг (Dynamic Inventory Allocation and Maintenance Optimizer, DINAMO) для управления ценами по всем направлениям. В People Express тоже использовались простые стратегии управления ценами для дифференциации сезонных и внесезонных тарифов, но их информационная система была намного проще, чем DINAMO, и отставала от нее по своей эффективности. В результате People Express начала терять до 50 миллионов долларов в месяц, что привело ее к банкротству и в конечном счете исчезновению в 1987 году, когда она была приобретена компанией Continental Airlines [Vasigh et al., 2013]. Однако American Airlines не только выиграла конкурентную борьбу с People Express, но и увеличила выручку на 14,5 % и прибыль на 47,8 % через год после внедрения системы DINAMO.

Случай с American Airlines стал первым крупным успехом в практике управления доходами. К началу 1990-х этот подход был взят на вооружение в других отраслях, где производятся скоропортящиеся товары или предоставляются услуги, заказываемые заранее: гостиничный бизнес, прокат автомобилей и даже

продажа телевизионной рекламы. Успех управления доходами в гражданской авиации явно связан с конкретными особенностями спроса и предложения в этой области:

- Спрос значительно различается между клиентами, рейсами и временем: покупательная способность бизнес-путешественников намного выше, чем у бюджетных путешественников, рейсы в сезон имеют бóльшую загрузку, чем в межсезонье, и т. д.
- Предложение, то есть свободные места, очень негибкий ресурс. Авиакомпании производят места большими блоками, планируя рейсы, и после того как рейс будет запланирован, число мест уже нельзя изменить. Непроданные места нельзя удалить, поэтому прибыль авиакомпании полностью зависит от ее способности эффективно управлять спросом и продажами.

Из вышесказанного можно сделать вывод, что управление доходами можно рассматривать как аналог управления цепочками поставок или, если хотите, цепочками управления спросом, целью которого является борьба с негибкостью производства и его адаптации к требованиям рынка (и наоборот, манипулирует спросом, чтобы привести его в соответствие с предложением), как показано на рис. 1.3.

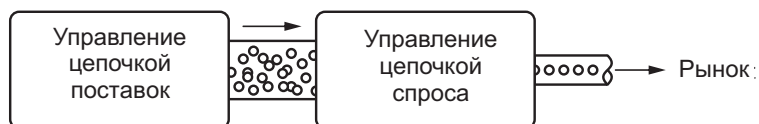


Рис. 1.3. Управление доходами как аналог управления цепочкой поставок

Этот конфликт между спросом и предложением можно увидеть не только в гражданской авиации, но и в других отраслях. Гостиницы и прокат автомобилей — самые близкие примеры, но реклама, ретейл и другие отрасли также демонстрируют особенности, указывающие на применимость алгоритмических методов для управления цепочкой спроса.

1.3.3. Наука маркетинга

Примеры онлайн-рекламы и продажи билетов в авиакомпаниях дают общее представление, насколько продвинулось применение алгоритмических методов в индустрии. Их внедрение подкреплялось быстрым развитием науки о маркетинге и, наоборот, развитие научных методов маркетинга подталкивалось и стимулировалось практическими потребностями. Маркетинг как дисциплина возник в начале 1900-х и в первые пять десятилетий своего существования ориентировался в основном на

описательном анализе процессов производства и распределения, то есть на сборе фактов о потоках товаров от производителей к потребителям. Идея о возможности подкрепления маркетинговых решений методами математического моделирования и оптимизации начала завоевывать популярность в 1960-х годах, что можно объяснить несколькими факторами. Во-первых, на науку о маркетинге повлияли достижения в области *исследования операций* — дисциплины, которая занимается проблемами выбора оптимальных решений в военной сфере и в бизнесе за счет применения статистического анализа и математической оптимизации. Исследование операций в свою очередь зародилось в годы Второй мировой войны в контексте планирования войсковых операций и оптимизации ресурсов. Во-вторых, развитие математических методов в маркетинге можно объяснить технологическими изменениями и внедрением первых мейнфреймов (мощных вычислительных систем) в организациях, позволивших собирать большие объемы данных и реализовать алгоритмы анализа и оптимизации. Наконец, специалисты по маркетингу начали осознавать, что старые методы продаж исчерпали себя и маркетинг необходимо переопределить как смешение ингредиентов, которые можно контролировать и оптимизировать; так в 1960 году появилось понятие маркетинг-микса. Наука о маркетинге бурно развивалась в шестидесятых и семидесятых годах XX века, когда были разработаны многочисленные количественные модели ценообразования, распределения и планирования продукции с использованием вероятностных и оптимизационных методов. Некоторые из этих методов с легкостью внедрялись в практику, как в случае с управлением доходами в авиакомпаниях и гостиничном бизнесе, но многие другие имели ограниченную практическую применимость, поэтому во многих отраслях общий уровень автоматизации оставался на довольно низком уровне [Wierenga, 2010].

Развитие цифровых каналов кардинально изменило ситуацию. Цифровые медиа обусловили необходимость и дали возможность принимать миллионы микро-решений на уровне отдельных клиентов и оказывать совершенно новые услуги, такие как поиск товаров или мобильные уведомления в режиме реального времени. Это породило проблемы, которые часто выходят за рамки экономического моделирования и оптимизации — основной задачи науки о маркетинге — и требуют использования передовых методов разработки программного обеспечения и анализа данных, изначально не связанных с маркетингом. В современной розничной торговле (ритейле), например, значительная часть выручки может быть получена за счет оказания услуг поиска и рекомендаций, которые внутренне опираются на методы анализа текста, а не на экономические модели. Многие такие методы пришли из сфер, далеких от маркетинга, таких как биология и исследование генома. Подводя итог, можно сказать, что традиционное экономическое моделирование, наука о данных, разработка ПО и классические приемы маркетинга очень важны для создания программных систем.

1.4. Программные услуги

Модель маркетинг-микса определяет четыре фактора, влияющих на решение потребителя о покупке, которые могут контролироваться компанией: продукт, продвижение, цена и распространение. Однако такая классификация слишком широка и дает мало указаний, как именно следует строить программную маркетинговую систему. На данный момент мы уже знаем, что программную систему можно рассматривать как поставщика одной или нескольких функциональных услуг, реализующих определенные бизнес-процессы, такие как ценообразование или продвижение. Следовательно, можно конкретизировать задачу, определив набор услуг, каждая из которых реализует определенную функцию и имеет свой вход (цели) и выход (действия). Существуют разные варианты деления маркетинг-микса на функциональные услуги, в зависимости от отраслевой принадлежности компании и ее бизнес-модели. Мы решили выделить шесть основных функциональных услуг, актуальных для широкого спектра вертикалей управления взаимоотношениями с потребителями (Business-to-Consumer, B2C): продвижение, реклама, поиск, рекомендации, цена и планирование ассортимента. Эти шесть услуг являются главной темой данной книги, и все следующие главы будут посвящены обсуждению вопросов их проектирования и создания. К разным услугам применяются разные принципы проектирования и реализации, но при этом они связаны множеством связей. Давайте кратко перечислим эти связи и познакомимся с некоторыми общими рекомендациями по проектированию, которые более подробно будут рассматриваться в остальной части книги.

Отношения между шестью услугами, которые мы определили, а также их связь с маркетинг-миксом можно обозначить так:

- Основная задача услуг продвижения и рекламы — оценить пожелания потребителя и донести до него правильную информацию. Для этого обычно требуется выявить потребителей, которых можно побудить выполнить определенные действия, способствующие достижению желаемой бизнес-цели. Способность идентификации подходящих потребителей и выбора правильных предложений для них является краеугольным камнем этой группы услуг. С точки зрения маркетинг-микса, эти услуги относятся непосредственно к сфере «Продвижение» и связаны со сферой «Цена» через затраты и прибыль, связанные с продвижением и рекламными кампаниями.
- Услуги поиска и рекомендаций решают задачу поиска нужных товаров для данного потребителя и являются естественным продолжением предыдущей группы услуг. Основная их цель — обеспечить и упростить поиск товара. В маркетинг-миксе они связаны со сферами «Распространение» и «Продвижение». Эта группа услуг нуждается в понимании намерений потенциального

покупателя, выражаемых явно или опосредованно, и в способности находить предложения, соответствующие этим намерениям.

- Целью услуг цены и планирования ассортимента является определение и оптимизация набора предложений и их свойств, включая цену. Эти услуги часто опираются на возможность прогнозирования спроса в зависимости от ассортимента, цен и других параметров, что дает возможность проводить анализ возможных альтернатив и оптимизировать разные варианты. В маркетинг-миксе эта группа главным образом охватывает сферы «Цена» и «Товар».

Такое деление на три группы, показанное на рис. 1.4, удобно тем, что отражает сходство целей услуг и принципов их проектирования. Структура остальной части книги подчинена этому делению, чтобы охватить как фундаментальные возможности, такие как правильное определение потребителей для выработки предложений, так и отдельные услуги.

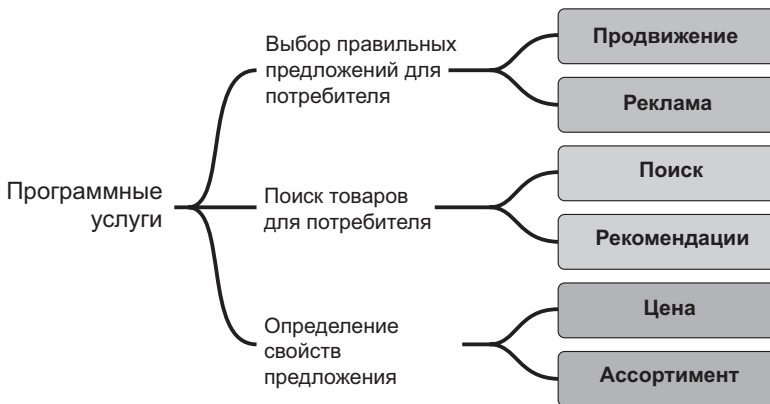


Рис. 1.4. Программные услуги

Следующий вопрос, который можно задать, — что общего между этими услугами с точки зрения проектирования и реализации. Несмотря на то что принципы проектирования и методы реализации сильно разнятся для разных услуг, программный подход вводит общие руководящие принципы, которым явно или косвенно могут следовать все услуги.

На основе этих принципов можно определить базовую терминологию и компоненты, которые в дальнейшем можно разработать в рамках соответствующих сфер.

Идея программного метода подчеркивает целенаправленный подход к проектированию, поэтому попробуем определить общую структуру, начав с понятия

бизнес-цели. Чтобы понять и выполнить цель, программная услуга, вероятно, должна включать определенный набор функциональных компонентов, которые могут иметь разную архитектуру для разных сфер (рис. 1.5):

- Поскольку для автоматического принятия решений необходимы данные, конвейер принятия решений начинается со сбора данных. В числе примеров исходных данных для большинства маркетинговых приложений можно назвать личные сведения о потребителях, профили их поведения, учетные данные и сведения о покупках.
- Исходные данные часто необходимо преобразовать в четко определенные *признаки*, которые затем можно передать на вход алгоритмов анализа и принятия решений. Это объясняется тем, что программные услуги часто полагаются на некоторые меры сходства между сущностями, такими как товары или потребители, для выявления закономерностей и принятия решений, что требует, чтобы сущности были представлены в виде сопоставимых наборов атрибутов. Услуги поиска, например, нередко опираются на некоторую меру сходства между запросом пользователя и товарами для выбора наиболее подходящих предложений, что в свою очередь требует преобразования запросов и информации о товарах в четко определенные, сопоставимые представления. Разработка этих атрибутов, которые во многих контекстах называют признаками, играет решающую роль в программных системах.
- Наиболее важным этапом в программном конвейере является оценка соответствия разных стратегий бизнес-целям. Как правило, для этого требуется разработать одну или несколько моделей, которые получают потенциальные решения и производят сигналы, определяющие уровень соответствия. Например, услуга продвижения может опираться на *модель*, оценивающую клиента с точки зрения его склонности к приобретению определенного товара, услуга ценообразования может оценить разные варианты цены в соответствии с ожидаемой прибылью, а услуга поиска может оценить релевантность товаров запросу.
- *Сигналы*, генерируемые моделями, несут информацию о качестве разных возможных решений. Однако для бизнес-операции часто требуется объединить множество сигналов и промежуточных решений в заключительный план действий. Например, маркетолог может проводить рекламные акции для наиболее ценных клиентов, что требует оценки отдельных потребителей, но окончательный список рассылки должен ограничиваться бюджетом компании. То есть программные услуги обычно содержат компонент оптимизации или смешивания сигналов, который принимает окончательные решения.
- Программная услуга взаимодействует с внешним миром через *маркетинговые каналы* и интегрируется с другими услугами. Эти каналы определяют

набор возможных бизнес-операций, которые может выполнять услуга, и параметры операций, которыми услуга может управлять, такие как уровни цен, величины скидок, сообщения электронной почты или порядок товаров в списке с результатами поиска. Эти возможности управления используются программной услугой для выполнения принятых решений, поэтому решения, в конечном итоге, должны выражаться в виде параметров, доступных для управления.

- Наконец, обратная связь, полученная по каналам исполнения, может направляться обратно в модели и процедуры оптимизации для учета результатов и корректировки логики принятия решений. Этап оценки является обязательной частью всех маркетинговых услуг, и многие маркетинговые подходы опираются на способ оптимизации методом проб и ошибок.

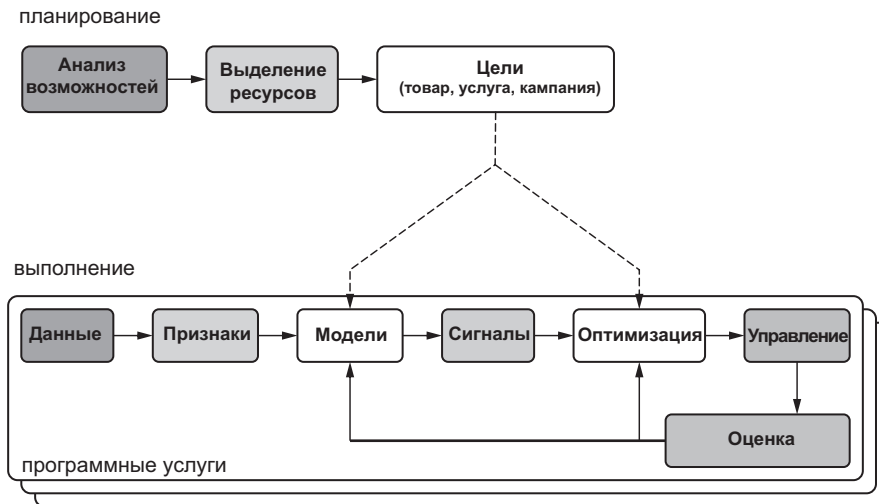


Рис. 1.5. Обобщенная структура программных услуг

На практике маркетолог часто стремится достичь нескольких стратегических целей и может использовать для этого несколько услуг. Программная инфраструктура может способствовать этому, предоставляя возможности для анализа конъюнктуры рынка и глобального распределения ресурсов, чтобы помочь выработать цели и параметры для отдельных услуг, а также консолидировать оценки. Эта возможность планирования и конвейеры исполнения отдельных услуг образуют полную программную экосистему. Инфраструктура, которую мы определили к данному моменту, пока выглядит весьма абстрактно, но в последующих главах мы постепенно определим методы проектирования для всех этапов конвейера.

1.5. Кому адресована эта книга?

Эта книга адресована всем, кто хочет научиться создавать продвинутые маркетинговые программные системы. Она будет полезна специалистам по маркетингу и разработчикам, однако писалась с прицелом на две целевые аудитории. Первая — разработчики маркетингового программного обеспечения, руководители проектов и программисты, желающие познакомиться с методиками, которые можно использовать в программных продуктах для маркетинга, а также узнать об экономических основах этих методик. Вторая — рыночные аналитики и руководители, нуждающиеся в рекомендациях, как маркетинговые организации и службы могут извлечь выгоду из машинного обучения и больших данных и как современные предприятия могут использовать передовые методы автоматизации принятия решений.

Предполагается, что читатели имеют как минимум начальный уровень подготовки в статистике, численных методах и программировании. Хотя большинство методов, представленных в книге, использует относительно простую математику, эта книга может не подойти читателям, которых интересуют только бизнес-аспекты маркетинга: это не традиционная книга о маркетинге — она рассказывает об *автоматизации маркетинга*. В целом, если следующая формула не вызывает у вас неприятия, вы наверняка справитесь с этой книгой:

$$\mathbb{E}[\chi] = \int_{-\infty}^{\infty} x f_{\chi}(x) dx.$$

Эта книга состоит из шести глав. В первой главе, «Введение», излагаются основные понятия и принципы алгоритмического маркетинга и рассматриваются несколько вдохновляющих примеров, которые иллюстрируют предпосылки и преимущества алгоритмического подхода. Вторая глава, «Обзор предиктивного моделирования», посвящена математическим основам алгоритмического маркетинга. Последние четыре главы охватывают четыре разные сферы маркетинга — продвижение и реклама, услуги поиска, рекомендации и ценообразование, — по одной в каждой главе. Эти четыре главы следуют той же алгоритмической методологии и, соответственно, имеют одинаковую структуру: каждая глава начинается с описания окружения с целью изучения ограничений и переменных, которые можно оптимизировать, затем следует рассмотрение бизнес-целей для точного определения задач оптимизации и, наконец, изучаются методы автоматизации принятия решений для разных задач и сценариев в конкретной сфере. Поскольку каждая из четырех основных глав посвящена своей сфере, все они относительно независимы, и читатели могут использовать эту книгу как справочник, читая только разделы, которые больше соответствуют их нуждам, или последовательно, от корки до корки.

Читатели могут свободно пропускать части, не соответствующие их подготовке или специализации. Например, читатели, знакомые с теорией вероятностей,

математической статистикой и машинным обучением могут бегло просмотреть вторую главу или вообще пропустить ее. Те, кого в первую очередь интересуют коммерческие аспекты и возможности алгоритмического подхода, могут прочитать разделы с описанием окружения, бизнес-целей и проблем оптимизации и пропустить математические детали с примерами вычислений. С другой стороны, читатели, занимающиеся реализацией маркетинговых систем, должны уделить особое внимание алгоритмам, вычислительным примерам и деталям реализации.

1.6. Итоги

- Программный подход к маркетингу фокусируется на создании высокоавтоматизированных систем и процессов, управляемых бизнес-целями. Программные методы могут применяться во всех сферах маркетинг-микса: продукт, продвижение, цена и распространение.
- Программный маркетинг можно рассматривать как набор услуг, позволяющий получить представление о состоянии рынка и взаимодействовать с ним. Эти услуги могут использоваться внутри компании, владеющей базой с информацией о клиентах, или быть проданы третьей стороне. Программные компоненты должны быть самодостаточными, комплексными и способными реализовать достаточно высокий уровень абстракции для продажи в качестве высокоценных услуг.
- Программная услуга часто выступает в роли динамического регулятора спроса и, следовательно, аналога управления цепочкой поставок. Эффективность программных методов возрастает с изменением спроса потребителей или со временем и уменьшается с увеличением гибкости производства. Следовательно, программные методы приносят наибольшую пользу в областях с большим разнообразием потребительских вкусов и доходов и/или в областях с негибким производством, где создается избыток товаров или услуг с высокими издержками при падении спроса.
- Наиболее важные примеры программных услуг включают продвижение, рекламу, поиск, рекомендации, ценообразование и планирование ассортимента. Принципы проектирования услуг отличаются, но некоторые функциональные возможности и логические компоненты присутствуют во всех услугах. Примерами таких компонентов являются модели количественной оценки, которые оценивают соответствие возможных решений выбранной бизнес-цели, модели оптимизации, анализирующие и смешивающие оценки для принятия окончательного решения, и компоненты управления, используемые для преобразования решений в действия.

2

Обзор предиктивного моделирования

Алгоритмический маркетинг по определению не может существовать без методологии оценки возможных бизнес-действий и их результатов на основе доступных данных. В этой главе мы рассмотрим основные методы машинного обучения и экономического моделирования, которые позволяют анализировать будущие тенденции и образуют фундамент для остальной части книги. Цель этой главы — описание основных возможностей и ограничений предиктивного моделирования, а не исчерпывающее исследование алгоритмов машинного обучения. Мы опишем только несколько методов, обычно используемых в маркетинговых приложениях, и представим математический аппарат и примеры только для иллюстрации возможностей, ограничений и взаимосвязей с другими методами. В этой главе мы не будем вдаваться в практические аспекты моделирования, такие как подготовка входных данных и оценка получившихся моделей; все эти тонкости более подробно освещаются в остальной части книги.

2.1. Описательная, предиктивная и предписывающая аналитика

Прежде чем перейти к предиктивному моделированию, кратко остановимся на терминологии, используемой в маркетинге. В бизнесе методы анализа данных часто подразделяются на три общих категории: описательные, предиктивные и предписывающие. К *описательным* относятся методы обобщения данных, оценки их качества и поиска корреляций. Примерами описательной аналитики являются управленческие отчеты, предоставляющие укрупненные данные о продажах, и методы анализа рыночной корзины, дающие информацию о товарах, часто приобретаемых вместе.

Описательная аналитика не преследует цели объяснить, как можно повлиять или оптимизировать наблюдаемые результаты. Задача *предиктивной аналитики* — оценить вероятность потенциального результата на основе наблюдаемых данных или известных до результата. Типичными примерами предиктивной аналитики могут служить прогнозирование спроса и оценка склонности для определения вероятной реакции потребителей на рекламную акцию. Обратите внимание, что слово «предиктивный» необязательно означает прогнозирование будущего — оно используется в том смысле, что предиктивная модель может оценить величину некоторого выходного параметра при некотором изменении входного параметра. Наконец, под *предписывающей аналитикой* подразумевается моделирование зависимостей между решениями и будущими результатами с целью оптимизации решений. Одним из основных примеров предписывающей аналитики является оптимизация цены, когда прибыль моделируется как функция цены, чтобы оценить, сколько долларов прибыли даст каждый доллар скидки и определить оптимальную величину скидки.

Действия и процессы, связанные с данными, в маркетинге обычно рассматриваются через призму этих трех видов аналитики, и все они играют важную роль. В программных приложениях, где ключом является автоматизация принятия решений, основное внимание уделяется предписывающей аналитике, которая в свою очередь основана на предиктивном моделировании. Поэтому важно понимать, что программные приложения используют только подмножество аналитических методов.

2.2. Экономическая оптимизация

Маркетинг — это деятельность, направленная на достижение определенных бизнес-целей за счет выполнения определенных бизнес-действий. Приступая к рассмотрению алгоритмического подхода, первое, что мы должны сделать, — перевести этот бизнес-язык в более формальные модели, описывающие конечную цель, пространство возможных действий и имеющиеся ограничения. Большинство маркетинговых задач, переведенные на этот эконометрический язык, естественным образом превращаются в задачи оптимизации, которые выражают бизнес-параметры (например, доход) как функцию возможных действий (например, маркетинговые кампании или корректировка ассортимента) и отыскивают оптимальное воздействие среди возможных стратегий.

Экономическая модель также является функцией данных, в том смысле, что она использует свойства и параметры, извлеченные из прошлого опыта. Например, представьте ретейлера, планирующего рекламную рассылку. Пространство возможных действий можно определить как набор решений (отправлять/не отправлять приглашение), принимаемых для отдельных клиентов. При этом доход компании

зависит как от действий, то есть от того, будет ли отправлено приглашение, так и от данных, таких как ожидаемый доход от данного клиента и почтовые расходы. Этот подход можно выразить более формально:

$$s_{opt} = \operatorname{argmax}_{s \in S} G(s, D), \quad (2.1)$$

где D — данные, доступные для анализа, S — пространство действий и решений, G — *экономическая модель*, отображающая действия и данные в экономический результат, и s_{opt} — оптимальная стратегия. Архитектура модели G полностью зависит от области применения. В следующих главах мы обсудим конструкцию различных моделей в контексте конкретных маркетинговых задач, однако есть несколько общих соображений, которые следует учитывать при проектировании любых моделей.

Во-первых, нужно определить *бизнес-цель* и выразить ее в виде числовой метрики, которая может быть предметом оптимизации. Во многих случаях целесообразно моделировать и оптимизировать прибыль, но, как будет показано позже, в некоторых случаях могут ставиться другие цели. Проектирование цели может оказаться непростой задачей, если она представляет компромисс между прибылью предприятия и полезностью для потребителя, как, например, поиск в онлайн-каталоге, который должен возвращать результаты, релевантные запросу потребителя и соответствующие целям и правилам продвижения товаров.

Во-вторых, необходимо учесть доступность данных или решить проблему их сбора. Роль сбора данных в уравнении 2.1 имеет большое значение, потому что модель G связей действий с результатами может быть очень сложной и определяться на основе данных с применением регрессионного анализа или других методов выявления скрытых закономерностей. В некоторых случаях модель невозможно определить полностью либо из-за высокой сложности (например, поведение пользователя нельзя точно предсказать), либо из-за невозможности экстраполировать имеющиеся данные на конкретный случай (например, когда действие заключается в представлении совершенно нового товара или услуги). В любом случае данные следует рассматривать как бизнес-актив, который может потребовать инвестиций и поиска компромисса между затратами на получение данных и ценностью этих данных. Например, параллельная проверка моделей в реальности — простой способ пожертвовать экономической эффективностью (одновременный запуск нескольких моделей, очевидно, не самое оптимальное решение), чтобы получить больший объем данных и, иногда, более простую модель.

В-третьих, модель можно создать с разными уровнями детализации. Одну и ту же цель можно выразить, используя разные модели, в зависимости от пространства возможных действий, доступных данных и знаний о бизнес-ограничениях. Одним из ключевых соображений при разработке модели является уровень укрупнения

данных. Классические экономические модели часто оперируют небольшим числом высокоуровневых сводных показателей, таких как общий спрос. Это упрощает модели с точки зрения вычислений и сбора данных, но ограничивает их способность моделировать сложные зависимости. Алгоритмические методы предполагают наличие мощной вычислительной инфраструктуры и данных высокого разрешения, обеспечивающих возможность более детального моделирования. Разницу между этими двумя подходами можно проиллюстрировать на следующем упрощенном примере [Kleinberg et al., 1998]. Допустим, что ретейлер продает товар с маржой m , а q_u — месячная сумма, вырученная за этот товар от потребителя u . Месячный доход в этом случае можно описать так:

$$G = \sum_u q_u m. \quad (2.2)$$

Ретейлер решил запустить рекламную кампанию со стоимостью одной акции c , чтобы увеличить продажи в k раз. Ретейлер может контролировать обе величины, k и c , выбирая более или менее агрессивную стратегию продвижения. Соответственно, решение задачи оптимизации можно определить так:

$$\max_s \sum_u k \cdot q_u m - c, \quad (2.3)$$

где s — стратегия продвижения, определяемая парой параметров k и c . Как видите, определение в терминах отдельных потребителей избыточно и данную задачу можно переопределить в терминах укрупненных параметров, а именно общего спроса

$$Q = \sum_u q_u \quad (2.4)$$

и общего бюджета кампании C . То есть задача определяется так:

$$\max_s k \cdot Q \cdot m - C. \quad (2.5)$$

Теперь предположим, что ретейлер хочет создать два разных потребительских сегмента и одному назначить стратегию $s_i = (k_i, c_i)$, а другому — стратегию $s_j = (k_j, c_j)$. Если допустить, что стратегии выбраны из пространства S , задачу оптимизации можно выразить так:

$$\max_{s_i, s_j \in S} \sum_u \max\{q_u k_i m - c_i, q_u k_j m - c_j\}. \quad (2.6)$$

Это выражение нелинейно в отношении параметров k и c , поэтому его непросто переопределить в терминах укрупненных параметров. Следовательно, может потребоваться использовать более сложные методы извлечения данных. Этот компромисс между укрупненными сводными данными и данными высокого разрешения

является общей закономерностью, возникающей во многих задачах из-за нелинейности зависимостей между переменными, вовлеченными в процесс оптимизации.

Наконец, следует отметить, что задача оптимизации 2.1 в целом некоторым образом зависит от времени из-за изменения окружения (на рынке появляются новые товары, конкуренты делают свои ходы и т. д.) и собственных действий предприятия. Один из возможных подходов к учету этой зависимости является использование модели без состояния, которую можно рассматривать как математическую функцию, принимающую аргументы, зависящие от времени, для получения эффектов памяти. Например, модель прогнозирования спроса может предсказывать спрос на следующий месяц, принимая аргументы с уровнями скидок за последнюю неделю, последние две недели, последние три недели и т. д.

2.3. Машинное обучение

В предыдущем разделе отмечалось, что цель оптимизации можно определить как функцию от данных и стратегии маркетинга, $G(s, D)$. Наш следующий шаг — дать более формальное определение данных, которое поможет преодолеть разрыв между экономической моделью и методами извлечения данных.

Первое, на что нужно обратить внимание: процесс экономического моделирования касается только определенных показателей компании или потребителя, непосредственно связанных с моделируемой целью. Примерами таких показателей являются спрос на определенный товар или склонность потребителя реагировать на рекламную акцию. В большинстве случаев маркетинговая стратегия и действия не определяют эти показатели напрямую, а только влияют на них. Скидка, например, может увеличить спрос на определенный товар, а может не увеличить, если аналогичную скидку одновременно решили предложить конкуренты. Следовательно, мы заинтересованы в поиске функциональной зависимости между контролируруемыми и неконтролируемыми факторами и интересующими нас параметрами. В терминах теории вероятностей это можно выразить как условное распределение:

$$p(y | \mathbf{x}), \quad (2.7)$$

где \mathbf{x} — вектор факторов, а y — показатели. В примере со скидками \mathbf{x} может включать такие переменные, как цена на данный товар, цены на сопутствующие товары и цены конкурентов, а y выражать спрос, измеряемый в проданных единицах. Каждая маркетинговая стратегия s в этом случае соответствует определенной комбинации факторов, то есть некоторому вектору \mathbf{x} , который мы обозначим как $\mathbf{x}(s)$. В предположении, что распределение известно, задачу экономической оптимизации можно переписать в терминах, введенных выше:

$$\max_s G(p(y|\mathbf{x}(s))). \quad (2.8)$$

Если условное распределение $p(y|\mathbf{x})$ оказывается сложным и должно извлекаться из данных, а не определяться вручную, в игру вступают данные. В этом случае нас интересуют данные с парами (\mathbf{x}, y) , которые извлекаются из истинного, но неизвестного распределения $p_{data}(y|\mathbf{x})$. Мы будем называть эти пары *образцами* или *точками данных*. Входные данные часто имеют форму, непригодную или неоптимальную для моделирования, поэтому вектор \mathbf{x} и показатель y обычно конструируются из входных данных с применением очищающих и нормализующих преобразований. Элементы такого подготовленного вектора \mathbf{x} мы будем называть *признаками* или *независимыми переменными*, а y — *меткой отклика* или *зависимой переменной*. В предположении наличия n точек данных и m признаков все векторы признаков можно представить в виде матрицы \mathbf{X} с размерами $n \times m$, которая называется *матрицей плана*, а все переменные отклика — в виде n -мерного вектора-столбца \mathbf{y} . Каждая строка в матрице плана \mathbf{X} — это вектор признаков \mathbf{x} , а каждый элемент \mathbf{y} — метка отклика y . Все точки данных можно представить в виде матрицы \mathbf{D} с размерами $n \times (m + 1)$, имеющей следующую структуру:

$$\mathbf{D} = [\mathbf{X} | \mathbf{y}] = \left[\begin{array}{ccc|c} - & x_1 & - & y_1 \\ - & x_2 & - & y_2 \\ & \vdots & & \vdots \\ - & x_n & - & y_n \end{array} \right]. \quad (2.9)$$

Наша цель — создать статистическую модель, аппроксимирующую истинное распределение $p_{data}(y|\mathbf{x})$ распределением $p_{model}(y|\mathbf{x})$, полученным из данных, то есть фактически экономическая модель будет представлена следующей аппроксимацией:

$$\max_s G(P_{model}(y|\mathbf{x}(s))). \quad (2.10)$$

Во многих практических применениях требуется указать не распределение, а оценить наиболее вероятное значение y на основе \mathbf{x} , то есть получить функцию

$$\hat{y} = y(\mathbf{x}), \quad (2.11)$$

где слева находится оценочное значение переменной отклика. Оценить наиболее вероятное значение y проще, чем целое распределение, и, как будет показано ниже, это можно сделать без точной оценки фактических значений вероятностей. Экономическую модель можно оценить как

$$\max_s G(y(\mathbf{x}(s))). \quad (2.12)$$

В общем случае модель G может использовать несколько моделей данных, полученных из одного или нескольких наборов данных. Этот подход делит исходную задачу моделирования на следующие меньшие задачи, которые можно исследовать по отдельности:

- Распределение $p(y | \mathbf{x})$ должно оцениваться на основе данных. Это стандартная задача машинного обучения, относящаяся к классу задач *обучения с учителем*.
- В некоторых случаях соответствующие значения \mathbf{x} и y явно не представлены в доступных данных. Однако иногда можно найти преобразование, отображающее исходные данные в новое представление, более подходящее для моделирования. Эта задача известна как *проектирование признаков*. В некоторых случаях проектирование признаков можно выполнить полуавтоматически, с использованием относительно простых методов. Например, в одних случаях удастся повысить точность модели за счет использования логарифмов входных значений. В других приходится использовать более продвинутые методы машинного обучения, чтобы отыскать подходящее представление. Эта задача называется *обучением представлений*.
- Наконец, должна быть определена экономическая модель, оценивающая бизнес-результаты распределения. Это экономическая задача, а не задача машинного обучения, поэтому мы подробно обсудим ее в последующих главах, посвященных конкретным маркетинговым проблемам. Однако существует ряд стандартных моделей для решения базовых задач, таких как прогнозирование потребительского выбора, о которых рассказывается далее в этой главе.

А теперь продолжим обзор каждой из трех областей. Далее мы рассмотрим инструментарий, доступный разработчикам прикладных программных систем, и определим базовую терминологию, которую будем использовать в оставшейся части книги. Здесь мы сосредоточимся на концептуальных проблемах и решениях, не углубляясь в детали алгоритмов и реализации, которые можно найти в специализированных книгах по машинному обучению, например, таких как Bishop, 2006, Murphy, 2012 и Zaki and Meira, 2014.

2.4. Обучение с учителем

Выше мы видели, что задачу моделирования частично можно свести к поиску аппроксимации распределения $p(y | \mathbf{x})$ на основе имеющихся точек данных \mathbf{x} и y . Функцию p , отображающую m -мерный вектор признаков в значения вероятностей, можно интерпретировать как функцию плотности вероятности для непрерывного y или функцию вероятности для дискретного y . Во многих случаях не требуется иметь все распределение, достаточно функции, предсказывающей

наиболее вероятный отклик y для входного вектора \mathbf{x} . Задача поиска таких аппроксимаций распределения или функций называется обучением с учителем, потому что данные, содержащие переменные отклика, «направляют» процесс обучения. Различают два основных типа задач обучения с учителем. Если переменная отклика определяет категорию, то есть y принадлежит некоторому конечному множеству классов, задача называется задачей *классификации*. Если переменная отклика изменяется в непрерывном диапазоне значений, задача называется задачей *регрессии*.

В этом разделе мы сначала обсудим подходы к проблеме оценки распределений и предиктивных моделей, а также связь между истинным распределением $p_{data}(y|\mathbf{x})$ и его аппроксимацией $p_{model}(y|\mathbf{x})$, а затем рассмотрим несколько примеров конструирования моделей.

2.4.1. Параметрические и непараметрические модели

Одним из основных аспектов при проектировании предиктивной модели является выбор между параметрическим и непараметрическим подходами. Параметрический подход предполагает, что распределение данных имеет определенную функциональную форму, определяемую фиксированным числом параметров, поэтому задачу аппроксимации распределения можно выразить как задачу *обучения модели*, то есть подбор таких параметров, чтобы модель распределения

$$p_{model}(y|\mathbf{x}, \theta) \quad (2.13)$$

оптимально соответствовала данным. Конечно же, условие оптимальности тоже должно быть формально определено. Непараметрический подход предполагает увеличение числа параметров с увеличением объема обучающих данных, и в некоторых методах каждая точка данных может рассматриваться как параметр. Одним из наиболее часто используемых непараметрических методов является алгоритм k -ближайших соседей (kNN). Его идея состоит в том, чтобы предсказать значение переменной отклика для вектора признаков \mathbf{x} , опираясь на переменные отклика в обучающей выборке, которые являются ближайшими соседями \mathbf{x} в пространстве признаков. В задаче классификации вероятность принадлежности переменной отклика классу c можно оценить как

$$\Pr(y=c|\mathbf{x}, k) = \frac{1}{k} \sum_{i \in N_k(\mathbf{x})} \mathbb{I}(y_i=c), \quad (2.14)$$

где k — параметр алгоритма, определяющий количество соседей, $N_k(\mathbf{x})$ — k ближайших соседних точек данных в обучающем наборе данных и \mathbb{I} — функция-индикатор, равная 1, если ее аргумент истинен, и 0 в противном случае. Соседей

входного вектора \mathbf{x} можно определить с использованием любой метрики векторного расстояния, например евклидова расстояния. Решение о классификации затем можно принять, выбрав наиболее вероятный класс:

$$y = \underset{c}{\operatorname{argmax}} \Pr(y = c | \mathbf{x}, k). \quad (2.15)$$

Этот процесс изображен на рис. 2.1. Модель регрессии определяется аналогично, путем аппроксимации отклика как среднего значения переменных отклика соседей.

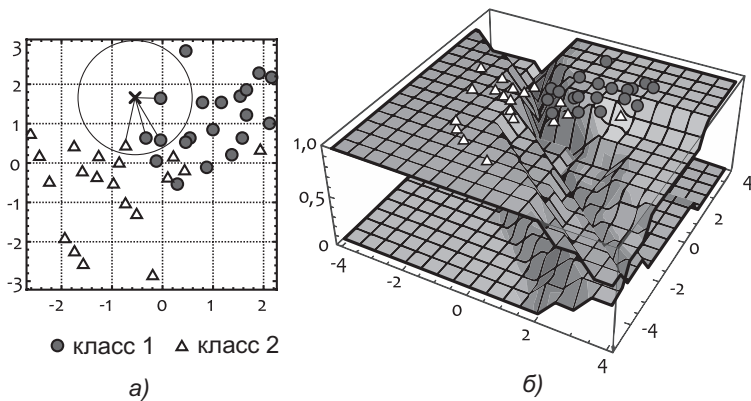


Рис. 2.1. Классификация двумерных точек с использованием алгоритма kNN.

а) Набор обучающих данных с двумя классами точек. Каждая точка в пространстве классифицируется по $k = 4$ ближайшим соседям. б) Вероятности классов, аппроксимируемые как функции от признаков, согласно уравнению 2.14

Алгоритм ближайших соседей является одним из самых простых в обучении с учителем; однако во многих ситуациях он может давать очень неплохие результаты. Например, он широко используется в алгоритмах выработки рекомендаций, о чем мы поговорим в соответствующей главе. Недостатком непараметрических методов является увеличение разреженности пространства с ростом размерности данных, из-за чего приходится рассматривать соседей, которые настолько далеки от данной точки, что на самом деле не способны надежно предсказывать зависимость между входами и выходами в требуемой области. Иначе говоря, модель k -ближайших соседей показывает хорошие результаты только для локальных наблюдений и не может обобщать закономерности, наблюдаемые в наборе данных. Эту проблему можно решить с помощью параметрических моделей, которые, хотя и имеют меньшую гибкость из-за ограниченного числа параметров, аппроксимируют свои параметры глобально.

2.4.2. Оценка методом максимального правдоподобия

Задача обучения модели сама по себе является задачей оптимизации, поэтому мы должны определить целевую функцию для оптимизации. Предположим, что набор из n точек данных извлекается независимо от распределения данных $p_{data}(\mathbf{x}, y)$. Каждый экземпляр данных (\mathbf{x}_i, y_i) можно интерпретировать как вход \mathbf{x}_i , дающий на выходе y_i . Целевую функцию в этом случае можно определить как вероятность наблюдаемого отклика, учитывая, что плотность вероятности, заданная параметрическим вектором θ , известна:

$$L(\theta) = P_{model}(\mathbf{y} | \mathbf{X}, \theta). \quad (2.16)$$

Эта функция называется *функцией правдоподобия*, или просто правдоподобием. Это вероятность наблюдения обучающих данных в предположении, что они получены из распределения, определяемого моделью с параметрами θ . Для анализа и расчетов часто удобнее использовать логарифм функции правдоподобия, известный как логарифм правдоподобия $LL(\theta)$. Наша цель состоит в том, чтобы найти вектор параметров, максимизирующий вероятность оценки:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p_{model}(\mathbf{y} | \mathbf{X}, \theta). \quad (2.17)$$

Предположив, что экземпляры независимы и подчиняются одному и тому же распределению, можно разбить вероятность правдоподобия на произведение n вероятностей отдельных точек данных:

$$LL(\theta) = \log p_{model}(\mathbf{y} | \mathbf{X}, \theta) = \sum_{i=1}^n \log p_{model}(y_i | \mathbf{x}_i, \theta). \quad (2.18)$$

Уравнение можно разделить на n , потому что оператор argmax безразличен к масштабированию. Далее θ_{ML} можно выразить в терминах математического ожидания для точек данных:

$$LL(\theta) = \mathbb{E}_{\mathbf{x}, y \sim p_{data}} [\log p_{model}(y | \mathbf{x}, \theta)]. \quad (2.19)$$

Теперь мы можем показать, что принцип максимального правдоподобия приводит к минимизации расхождения между распределением данных и его аппроксимацией. Стандартной мерой расхождения между двумя распределениями является расхождение Кульбака—Лейблера (Kullback—Leibler divergence), или KL -расхождение, которое определяется как

$$KL(p_{data}, p_{model}) = \mathbb{E}_{\mathbf{x}, y \sim p_{data}} \left[\log \frac{p_{data}(y | \mathbf{x})}{p_{model}(y | \mathbf{x}, \theta)} \right]. \quad (2.20)$$

Поскольку p_{data} не зависит от θ и не может быть предметом оптимизации, для минимизации расхождения достаточно минимизировать второй член, что эквивалентно максимизации логарифма правдоподобия из уравнения 2.19:

$$\operatorname{argmin}_{\theta} KL(p_{data}, p_{model}) = \operatorname{argmin}_{\theta} -LL(\theta). \quad (2.21)$$

То есть максимальное правдоподобие можно рассматривать как оптимизацию параметров модели для достижения соответствия распределения модели эмпирическому распределению.

2.4.3. Линейные модели

Принцип максимального правдоподобия обеспечивает общую основу для получения прикладных алгоритмов создания моделей данных. Теперь посмотрим, как этот принцип можно использовать для построения нескольких простых, но очень полезных предиктивных моделей. Сначала рассмотрим регрессионную задачу и реализуем модель, предсказывающую ответ y как непрерывную линейную функцию входного x . Затем перейдем к задаче классификации и обсудим несколько моделей, предсказывающих категорию подгонкой гиперплоскости, делящей пространство признаков на области, каждая из которых соответствует некоторому классу откликов.

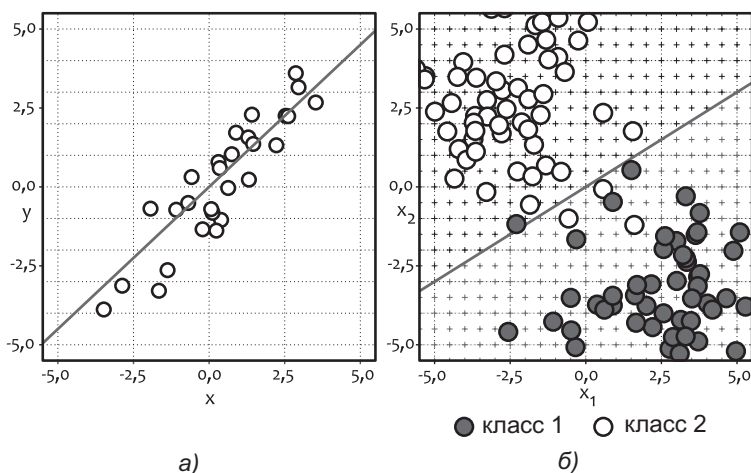


Рис. 2.2. а) Пример линейной регрессии в одномерном пространстве признаков.
б) Пример классификации с линейной границей решения и двумя двумерными пространствами признаков. Обучающие точки данных изображены окружностями

Все модели, которые мы рассмотрим, являются линейными. Регрессионная модель является линейной в том смысле, что зависимость между признаками и откликом моделируется линейной функцией, как показано на рис. 2.2. Если наблюдаемая зависимость на самом деле не является линейной, модель может неточно аппроксимировать данные. Модели классификации являются линейными в том смысле, что граница между классами моделируется как гиперплоскость, то есть данные *линейно неразделимы*. Если группы точек нельзя точно разделить гиперплоскостью, модели могут неправильно прогнозировать данные.

2.4.3.1. Линейная регрессия

Цель регрессионной модели — отобразить непрерывный вход \mathbf{x} в непрерывный выход y . В моделях линейной регрессии это отображение достигается за счет линейной функции от входных данных:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}, \quad (2.22)$$

где \mathbf{w} — вектор параметров модели, которые требуется получить. Таким образом, ошибка аппроксимации будет

$$\epsilon = y - y(\mathbf{x}) = y - \mathbf{w}^T \mathbf{x}. \quad (2.23)$$

Если допустить, что ошибка имеет нормальное распределение, тогда распределение оценки, произведенной моделью, также будет описываться законом нормального, гауссова, распределения со средним $\mathbf{w}^T \mathbf{x}$ и дисперсией σ^2 :

$$\begin{aligned} p(y | \mathbf{x}, \mathbf{w}) &= N(y | \mathbf{w}^T \mathbf{x}, \sigma^2) = \\ &= \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2} (y - \mathbf{w}^T \mathbf{x})^2 \right). \end{aligned} \quad (2.24)$$

Вставив такое распределение вероятности в определение логарифма правдоподобия (уравнение 2.18) и выполнив некоторые алгебраические преобразования, получим следующее выражение логарифма правдоподобия для максимизации:

$$\begin{aligned} LL(\mathbf{w}) &= \sum_{i=1}^n \log p(y | \mathbf{x}_i, \mathbf{w}) = \\ &= \sum_{i=1}^n \log N(y | \mathbf{w}^T \mathbf{x}_i, \sigma^2) = \\ &= -\frac{1}{2\sigma^2} SSE(\mathbf{w}) - \frac{n}{2} \log(2\pi\sigma^2), \end{aligned} \quad (2.25)$$

где SSE — сумма квадратов ошибок, которая определяется как

$$SSE(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2. \quad (2.26)$$

Интуитивно понятно, что максимизация правдоподобия эквивалентна минимизации ошибки. Предположив, что дисперсия фиксирована некоторым постоянным значением σ^2 (в общем случае ее тоже можно оценить), мы можем отбросить второй член в уравнении логарифмического правдоподобия 2.25, и тогда максимизацию логарифмического правдоподобия можно выполнить только относительно SSE . Сначала перепишем формулу правдоподобия в более компактной векторной форме:

$$\begin{aligned} LL(\mathbf{w}) &= -\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \\ &= \mathbf{w}^T (\mathbf{X}^T \mathbf{y}) - \frac{1}{2} \mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - \frac{1}{2} \mathbf{y}^T \mathbf{y}, \end{aligned} \quad (2.27)$$

где

$$\mathbf{X}^T \mathbf{X} = \sum_{i=1}^n x_i x_i^2 \quad \text{и} \quad \mathbf{X}^T \mathbf{y} = \sum_{i=1}^n x_i y_i. \quad (2.28)$$

Возьмем градиент относительно \mathbf{w} :

$$\nabla_{\mathbf{w}} LL(\mathbf{w}) = \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \mathbf{w}. \quad (2.29)$$

Приравняв градиент к нулю и решив уравнение относительно \mathbf{w} , мы получим аппроксимацию ML-оптимума \mathbf{w} в замкнутой форме:

$$\mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (2.30)$$

Линейная регрессия — самый простой метод предиктивного моделирования, но последовательность ее вывода является хорошей иллюстрацией, как можно выполнить обучение модели на основе принципа максимального правдоподобия. Позже мы увидим, что методы, способные аппроксимировать нелинейные зависимости, можно получить из линейной регрессии.

2.4.3.2. Логистическая регрессия и бинарная классификация

Второй пример, который мы рассмотрим, демонстрирует, как можно использовать принцип максимального правдоподобия для построения модели бинарной классификации, то есть модели, отображающей входной сигнал \mathbf{x} в один из двух возможных классов $y \in \{0, 1\}$. Пойдем по тому же пути, что и в случае с линейной регрессией, и определим форму модели. Наша цель — найти линейную границу решения (гиперплоскость), разделяющую два класса в точке, где

$$\Pr(y = 0 | \mathbf{x}) = \Pr(y = 1 | \mathbf{x}). \quad (2.31)$$

Это уравнение можно переписать в логарифмической форме:

$$\log \frac{\Pr(y = 0 | \mathbf{x})}{\Pr(y = 1 | \mathbf{x})} = 0. \quad (2.32)$$

Так как мы ищем линейную границу, гиперплоскость можно описать линейной функцией от входа \mathbf{x} с вектором коэффициентов \mathbf{w} :

$$\log \frac{\Pr(y = 0 | \mathbf{x})}{\Pr(y = 1 | \mathbf{x})} = \mathbf{w}^T \mathbf{x}. \quad (2.33)$$

Это означает, что вход можно отнести к классу 0, если $\mathbf{w}^T \mathbf{x}$ имеет положительное значение, и к классу 1, если $\mathbf{w}^T \mathbf{x}$ имеет отрицательное значение. Уравнение 2.33 эквивалентно выражению

$$\begin{aligned} \Pr(y = 1 | \mathbf{x}) &= \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \equiv g(\mathbf{w}^T \mathbf{x}), \\ \Pr(y = 0 | \mathbf{x}) &= \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})} = 1 - g(\mathbf{w}^T \mathbf{x}), \end{aligned} \quad (2.34)$$

где g — *логистическая функция*. Эта модель называется логистической регрессией. Обратите внимание, что это модель *классификации*, несмотря на название, сбивающее с толку. Далее мы должны вычислить логарифмическое правдоподобие для этого распределения:

$$\begin{aligned} LL(\mathbf{w}) &= \sum_{i=1}^n \log p(y_i | \mathbf{x}_i) = \\ &= \sum_{i=1}^n \log g(\mathbf{w}^T \mathbf{x}_i)^{y_i} (1 - g(\mathbf{w}^T \mathbf{x}_i))^{1-y_i} = \\ &= \sum_{i=1}^n y_i \log g(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \log (1 - g(\mathbf{w}^T \mathbf{x}_i)). \end{aligned} \quad (2.35)$$

Мы можем вычислить градиент этого выражения, но, к сожалению, не сможем получить оптимальное решение для \mathbf{w} в замкнутой форме, приравняв этот градиент к нулю и решая уравнение относительно \mathbf{w} . Поэтому для максимизации уравнения логарифмического правдоподобия 2.35 должны использоваться численные методы, такие как градиентный спуск, и аппроксимация оптимального веса \mathbf{w}_{ML} .

Логистическая регрессия моделирует вероятности классов с помощью логистической функции, определяемой уравнениями 2.34. Эта функция имеет S-образный график и также известна как сигмоидная кривая, крутизна которой определяется параметром ω . На рис. 2.3 показаны примеры логистических функций, где данные аппроксимируются двумя кривыми. Обратите внимание, что граница между классами — пересечение двух плоскостей — является прямой линией. Для этой иллюстрации использовались те же данные, что и в предыдущем примере с классификатором kNN на рис. 2.1, поэтому интересно сравнить поверхности вероятностей. Поверхности классификатора kNN имеют более сложную форму, чем поверхности логистической регрессии, потому что алгоритм kNN относится к непараметрическим методам. Поверхности, полученные в результате логистической регрессии, имеют гораздо более простую форму, определяемую логистической кривой.

Логистическая регрессия — один из простейших методов классификации и не способен аппроксимировать нелинейные границы между классами. Однако, как вы увидите далее, его и другие линейные методы можно расширить, добавив возможность аппроксимации нелинейных границ решений.

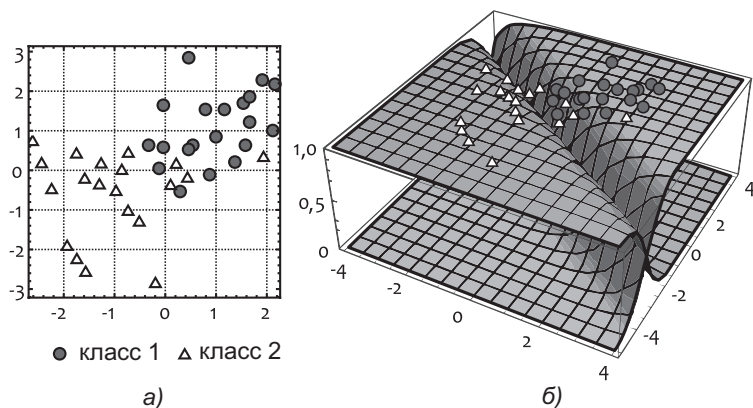


Рис. 2.3. Классификация двумерных точек с использованием логистической регрессии.
а) Обучающий набор данных с точками двух классов. б) Вероятности классов, аппроксимируемые как функции признаков, согласно уравнению 2.34

2.4.3.3. Логистическая регрессия и классификация с множеством категорий

Алгоритм логистической регрессии легко расширить для случая классификации с множеством категорий. Однако в этом случае, поскольку число классов больше двух, мы не можем использовать удобное отношение

$$\Pr(y = 0 | \mathbf{x}) = 1 - \Pr(y = 1 | \mathbf{x}), \quad (2.36)$$

применявшееся в примере, а одной линейной границы здесь недостаточно. Вместо этого можно оценить вероятность принадлежности к каждому классу c отдельно, используя выделенный вектор коэффициентов \mathbf{w}_c . То есть уравнение 2.33 можно переписать так:

$$\begin{aligned} \log \Pr(y = 0 | \mathbf{x}) &= \mathbf{w}_0^T \mathbf{x} - \log Z, \\ \log \Pr(y = 1 | \mathbf{x}) &= \mathbf{w}_1^T \mathbf{x} - \log Z, \\ &\vdots \\ \log \Pr(y = c | \mathbf{x}) &= \mathbf{w}_c^T \mathbf{x} - \log Z, \end{aligned} \quad (2.37)$$

где $\log Z$ — нормализующий член, гарантирующий, что полученное распределение y в действительности является распределением вероятностей с суммой, равной единице. Роль Z становится более четкой, если уравнение 2.37 записать в экспоненциальной форме:

$$\Pr(y = c | \mathbf{x}) = \frac{1}{Z} \exp(\mathbf{w}_c^T \mathbf{x}). \quad (2.38)$$

Таким образом, нормализующий множитель масштабирует распределение вероятностей, учитывая то обстоятельство, что сумма вероятностей всех классов должна быть равна единице:

$$\sum_c \Pr(y = c | \mathbf{x}) = 1. \quad (2.39)$$

Нормализующий множитель можно найти подстановкой 2.38 в 2.39:

$$\mathbf{Z} = \sum_c \exp(\mathbf{w}_c^T \mathbf{x}). \quad (2.40)$$

Подставив этот результат в уравнение 2.38, мы получим формулу оценки вероятности принадлежности к классу:

$$\Pr(y = c | \mathbf{x}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_i \exp(\mathbf{w}_i^T \mathbf{x})}. \quad (2.41)$$

Эту вероятность можно использовать для определения логарифмического правдоподобия, а векторы коэффициентов \mathbf{w}_i получить с применением численных методов, таких как градиентный спуск.

Важно отметить, что получившееся уравнение 2.41 можно интерпретировать как общий метод отображения вектора реальных значений в вектор вероятностей клас-

сов, независимо от базовой модели, производящей эти значения. Чтобы убедиться в этом, предположим, что некоторая модель, необязательно линейная, порождает вектор значений v , каждое значение в котором можно интерпретировать как относительный вес соответствующего класса. Эти веса необязательно нормализуются как вероятности: значения v могут выходить за диапазон $[0, 1]$ или их сумма может отличаться от единицы. Чтобы отобразить эти веса в нормализованные вероятности классов, определим обобщенную функцию, известную как функция *softmax*, которая повторяет уравнение 2.41, но использует вектор реальных значений в качестве параметра:

$$\text{softmax}(i, \mathbf{v}) = \frac{\exp(v_i)}{\sum_j \exp(v_j)}. \quad (2.42)$$

Нормализованные вероятности классов можно получить путем передачи весов классов v в функцию *softmax*. Мы будем использовать это свойство в последующих главах для конструирования предиктивных моделей.

2.4.3.4. Наивный байесовский классификатор

Последняя модель, которую мы рассмотрим в этом разделе, — наивный байесовский классификатор. Этот метод широко используется для классификации текста, поэтому его с успехом можно использовать при реализации услуг поиска и рекомендаций. Напомню, что задачу классификации с несколькими категориями можно определить как

$$\hat{y} = \underset{c}{\operatorname{argmax}} \Pr(y = c \mid \mathbf{x}). \quad (2.43)$$

Применяя правило Байеса к условной вероятности класса c , получаем следующее уравнение:

$$\Pr(y = c \mid \mathbf{x}) = \frac{\Pr(\mathbf{x} \mid y = c) \Pr(y = c)}{\Pr(\mathbf{x})}. \quad (2.44)$$

Вероятность вектора признаков в знаменателе одинакова для всех классов, поэтому ее можно отбросить:

$$\hat{y} = \underset{c}{\operatorname{argmax}} \Pr(y = c) \Pr(\mathbf{x} \mid y = c). \quad (2.45)$$

Важной особенностью наивного байесовского классификатора является предположение, что каждый признак x_i условно независим от любых других признаков

данного класса c . Это означает, что вероятность наблюдения вектора признаков \mathbf{x} для данного класса c можно разложить на множители:

$$\Pr(\mathbf{x} | y = c) = \prod_{i=1}^m \Pr(\mathbf{x}_i | y = c), \quad (2.46)$$

где m — длина вектора признаков. Это предположение, называемое *предположением условной независимости*, редко бывает верным на практике, тем не менее наивный байесовский алгоритм достаточно хорошо зарекомендовал себя в широком диапазоне подобных случаев. Например, этот классификатор с успехом применяется для классификации текста, когда каждый признак является словом, даже при том что слова в тексте не являются независимыми. Это можно объяснить тем, что наивный байесовский алгоритм продолжает оставаться корректным, даже если признаки взаимозависимы, но эти зависимости имеют определенную структуру и взаимно отменяют друг друга [Zhang, 2004]. Предположение о независимости позволяет переписать задачу классификации как

$$\hat{y} = \operatorname{argmax}_c \Pr(y = c) \prod_{i=1}^m \Pr(\mathbf{x}_i | y = c). \quad (2.47)$$

Параметрами данной модели являются значения $\Pr(y = c)$ и $\Pr(\mathbf{x}_i | y = c)$. Один из возможных подходов к обучению модели — рассматривать эти значения как неизвестные переменные, а не вероятности, и максимизировать логарифмическое правдоподобие, соответствующее уравнению 2.47. Однако легко показать, что это приводит к тем же результатам, что и в случае интерпретации параметров как эмпирических вероятностей. Другими словами, оценка максимального правдоподобия \hat{y} в 2.47 можно получить, если $\Pr(y = c)$ аппроксимировать как частоту класса c в обучающем наборе данных, а $\Pr(\mathbf{x}_i | y = c)$ — как частоту точек данных, принадлежащих классу c и имеющих значение признака \mathbf{x}_i . Это упрощает обучение наивной байесовской модели на практике.

В общем случае наивный байесовский классификатор не является линейным, но при определенных допущениях, верных для многих приложений, его можно считать линейным, поэтому часто он описывается как линейный. Например, рассмотрим случай допущения многомерности распределения $\Pr(\mathbf{x} | y = c)$. Оно верно, например, для случая классификации текста, когда каждый элемент вектора признаков является счетчиком слов. Вероятность для вектора признаков втекает из многомерного распределения с параметрическим вектором \mathbf{q}_c :

$$\Pr(\mathbf{x} | y = c) \propto \prod_{i=1}^m q_{ci}^{x_i}, \quad (2.48)$$

где q_{ci} — вероятность принадлежности значения признака x_i к классу c . Это выражение можно переписать в векторной форме:

$$\log \Pr(\mathbf{x} | y = c) = \mathbf{x}^T \log \mathbf{q}_c + \text{constant}. \quad (2.49)$$

Теперь покажем, что граница принятия решений между классами линейна, рассмотрим отношение вероятностей классов по аналогии с методом, который мы использовали для логистической регрессии. Предположив для упрощения, что имеются только два класса $y \in \{0, 1\}$, можно написать

$$\begin{aligned} \log \frac{\Pr(y = 1 | \mathbf{x})}{\Pr(y = 0 | \mathbf{x})} &= \log \Pr(y = 1 | \mathbf{x}) - \log \Pr(y = 0 | \mathbf{x}) = \\ &= \mathbf{x}^T (\log \mathbf{q}_1 - \log \mathbf{q}_0) + \log \Pr(y = 1) - \log \Pr(y = 0). \end{aligned} \quad (2.50)$$

Это линейная функция от \mathbf{x} , что доказывает линейность границы решения между классами.

2.4.4. Нелинейные модели

Линейные методы вполне подходят для многих маркетинговых приложений, и ответственность этих методов не следует недооценивать, но они плохо работают с наборами данных, имеющими нелинейные зависимости. Поэтому нам необходимо разработать методы, способные моделировать более сложные распределения. К этой задаче можно подойти с разных сторон, и здесь мы обсудим два основных семейства методов, которые часто используются на практике. Далее в книге мы обсудим еще несколько методов, таких как нейронные сети, в контексте конкретных областей применения.

2.4.4.1. Отображение признаков и ядерные методы

Линейность или нелинейность решаемых нами задач регрессии или классификации отражает характер моделируемого процесса, но вы должны понимать, что представление данных тоже может вносить свой вклад. Например, два линейно зависимых значения не были бы линейно зависимыми, если бы одно из них измерялось в логарифмической шкале. Верно и обратное утверждение: наборы данных, непригодные для обработки линейными методами, могут стать таковыми после отображения в другое пространство. Рассмотрим пример на рис. 2.4: одномерный набор данных слева состоит из двух классов, которые не являются линейно разделимыми, но двумерный набор данных, полученный из первого с помощью отображения $(x) \rightarrow (x, x^2)$, является линейно разделимым.

Такое преобразование исходного пространства признаков в другое, обычно более высокой размерности, называется *отображением признаков*. Интуитивно понятно, что добавление размерностей, определенных как нелинейные функции одного

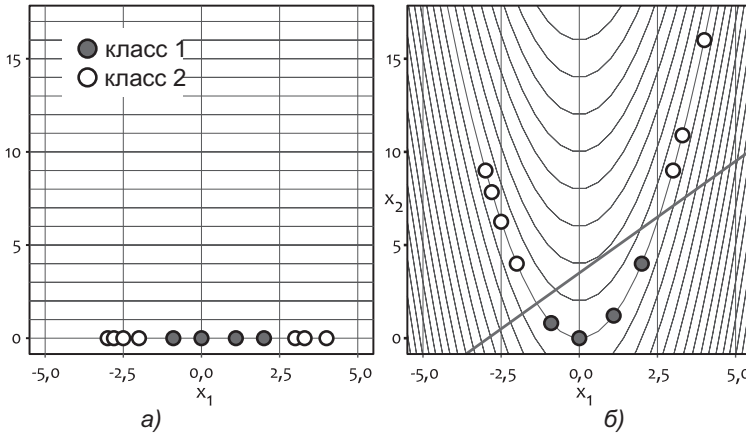


Рис. 2.4. Отображение признаков с использованием квадратичной функции.
 а) Исходные точки данных. б) Отображенные точки данных

или нескольких существующих признаков, обеспечивает большую гибкость для алгоритма регрессии или классификации, который мы пытаемся улучшить. Однако нам нужен метод определения функции отображения $\phi(\mathbf{x})$, которая создает новый вектор признаков более высокой размерности.

Первое, что мы должны отметить, — многие методы регрессии и классификации можно выразить в терминах расстояний между входным и обучающими векторами. Например, мы показали, что коэффициенты линейной регрессии можно вычислить как

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (2.51)$$

Это выражение можно умножить на единичную матрицу \mathbf{I} :

$$\mathbf{I} = \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad (2.52)$$

и, выполнив некоторые алгебраические преобразования, мы получим

$$\begin{aligned} \mathbf{w} &= \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \cdot (\mathbf{X}^T \mathbf{X}) \mathbf{X}^T \mathbf{y} = \\ &= \mathbf{X}^T \cdot \mathbf{a} = \sum_{i=1}^n a_i x_i, \end{aligned} \quad (2.53)$$

где вектор \mathbf{a} определяется как

$$\mathbf{a} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-2} \mathbf{X}^T \mathbf{y}. \quad (2.54)$$

Это означает, что переменную отклика можно аппроксимировать, используя только скалярные произведения между входным \mathbf{x} и обучающими \mathbf{x}_i векторами:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^n a_i \cdot \mathbf{x}^T \mathbf{x}_i. \quad (2.55)$$

Это очень важный результат, потому что теперь можно избежать явного отображения признаков и заменить его *ядерной функцией* (kernel function), инкапсулирующей вычисление скалярного произведения в отображаемом пространстве:

$$k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z}), \quad (2.56)$$

где \mathbf{x} и \mathbf{z} — два вектора признаков. Теперь уравнение 2.55 можно переписать исключительно в терминах ядерной функции:

$$y(\mathbf{x}) = \sum_{i=1}^n a_i k(\mathbf{x}, \mathbf{x}_i). \quad (2.57)$$

Проще говоря, мы заменили скалярное произведение в 2.55 функцией расстояния между векторами признаков. Теперь входные векторы можно не отображать, но, чтобы использовать ядерные функции вместо скалярного произведения, алгоритм необходимо модифицировать. Связь между ядерной функцией и функцией отображения можно проиллюстрировать на примере квадратичного ядра:

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T, \mathbf{z})^2. \quad (2.58)$$

Если исходное пространство признаков является двумерным, ядро расширяет его до трехмерного, содержащего как производные отдельных признаков, так и векторные произведения, которые могут сделать набор данных пригодным для обработки линейными методами:

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 = \\ &= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 = \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2) (z_1^2, \sqrt{2}z_1 z_2, z_2^2) = \\ &= \phi(\mathbf{x})^T \phi(\mathbf{z}), \end{aligned} \quad (2.59)$$

где

$$\phi(\mathbf{x}) = \{x_1^2, \sqrt{2}x_1 x_2, x_2^2\}. \quad (2.60)$$

Обратите внимание, что роль ядра играет простая функция расстояния между исходными векторами признаков, то есть лежащее в основе расширение размер-

ности полностью скрыто; благодаря этому можно разрабатывать ядра, которые соответствуют $\varphi(\mathbf{x})$ с очень большим или бесконечным числом размерностей, но остаются вычислительно простыми. Этот прием известен как *ядерный трюк* (kernel trick) и может использоваться для расширения многих методов машинного обучения. Выбор правильного ядра может быть непростой задачей, но есть несколько ядерных функций, которые, как известно, довольно универсальны и широко используются на практике. Выбор ядерной функции также зависит от приложения, потому что, по существу, эта функция описывает меру сходства между векторами признаков — ядра, хорошо работающие с профилями потребителей, могут не подойти для обработки текстовых описаний товаров и т. д.

Одним из самых известных членов семейства ядерных методов является метод опорных векторов (Support Vector Machines, SVM). В числе простейших алгоритмов SVM можно назвать методы линейной классификации и регрессии, но их легко привести к ядерной форме для аппроксимации нелинейных зависимостей. Рассмотрим пример классификатора SVM, изображенный на рис. 2.5. Он использует те же данные, что использовались выше в примерах классификации методом ближайших соседей и логистической регрессии, но имеет явно нелинейную границу решения.

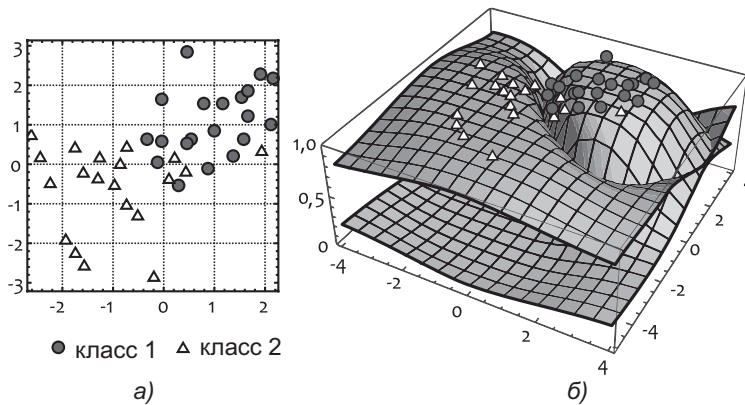


Рис. 2.5. Классификация двумерных точек с использованием метода опорных векторов и нелинейного ядра. а) Обучающий набор данных с двумя классами точек.
б) Вероятности классов как функция признаков

2.4.4.2. Адаптивный базис и деревья решений

Недостаток ядерного метода состоит в необходимости передавать ядерную функцию в виде параметра — ее нельзя получить в ходе обучения. Ядерная функция накладывает ограничения на форму границы решения, и хотя на практике можно попробовать несколько ядер или параметров ядра, чтобы найти хорошую аппрок-

симацию, иногда лучше выбрать другой подход. Мы можем сформулировать проблему обучения набора q базисных функций $\phi(\mathbf{x})$ так, чтобы ответ можно было предсказать как их взвешенную комбинацию:

$$y(\mathbf{x}) = \sum_{i=1}^q w_i \phi_i(\mathbf{x}). \quad (2.61)$$

При высокоадаптивных и нелинейных базисных функциях можно ожидать превосходства этого решения перед линейными методами в соответствующих задачах. Одной из наиболее широко используемых реализаций идеи адаптивного базиса являются деревья классификации и регрессии — семейство методов, генерирующих адаптивный базис $\phi_i(\mathbf{x})$ с использованием жадного эвристического алгоритма. Рассмотрим самый простой вариант этого решения.

Дерево классификации или регрессии создается рекурсивным разбиением пространства признаков на две части с использованием линейной границы решения, как показано на рис. 2.6. На каждом шаге рекурсии граница решения выбирается следующим образом:

- Сначала перечисляются гиперплоскости, кандидаты на роль границы. Один из возможных способов — опробовать все размерности (например, разбиение можно выполнить по горизонтали или по вертикали, как показано на рис. 2.6) и для каждого попробовать координаты всех точек данных в обучающем наборе.
- Граница-кандидат создает две области, каждая из которых может быть отмечена наиболее частым классом примеров в этой области или, в случае регрессии, средним значением переменных отклика. Затем эта метка используется в качестве прогнозируемого значения для любой точки данных, попадающей в область.
- Метка области используется для оценки качества границы-кандидата по ошибке классификации (соотношение между числом неправильно и правильно классифицированных примеров в области) или с применением другой метрики. Граница-кандидат выбирается по наибольшему количеству баллов.

После выбора границы алгоритм рекурсивно применяется к областям по обе ее стороны. Этот алгоритм создает набор отмеченных прямоугольных областей R_i , которые можно рассматривать как адаптивный базис. Чтобы убедиться в этом, можно переписать уравнение 2.61 с точки зрения регионов областей:

$$y(\mathbf{x}) = \sum_{i=1}^q w_i \mathbb{I}(\mathbf{x} \in R_i), \quad (2.62)$$

где w_i — метка области, а \mathbb{I} — 1 или 0, в зависимости от того, попадает ли \mathbf{x} в соответствующую область или нет. Деревья решений и более сложные производные

этого подхода, такие как *случайные леса*, обеспечивают мощное и широко используемое решение регрессии и классификации.

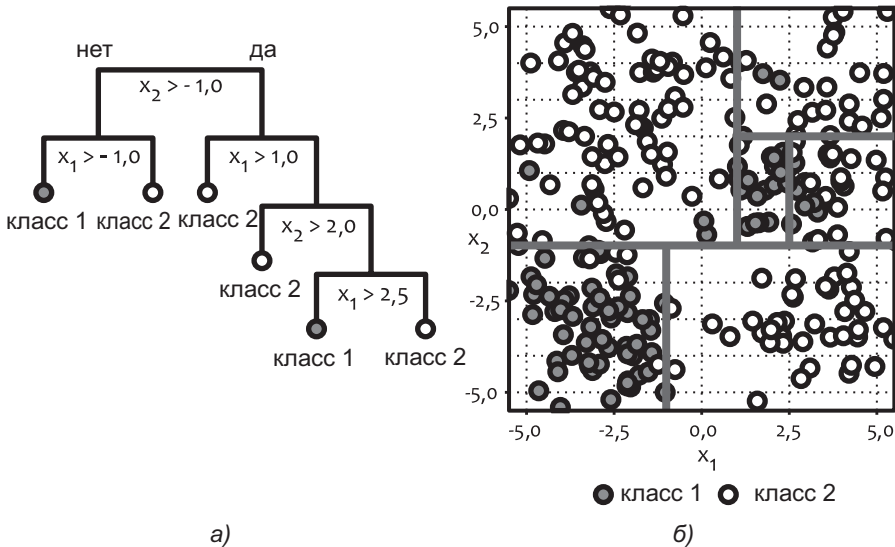


Рис. 2.6. Пример дерева классификации. а) Модель дерева классификации.
б) Обучающие данные и границы решения, соответствующие дереву

2.5. Обучение представлением

Методы обучения с учителем, которые мы только что рассмотрели, помогают описать зависимость между независимыми входными переменными и переменной отклика. Однако входные переменные могут быть избыточными и иметь запутанную структуру, усложняя исследование данных и обучение модели. Мы можем попытаться найти другое представление данных, которое лучше подходит для моделирования, удалив избыточности и корреляции, то есть распутав исходные данные.

Методы машинного обучения обычно подразделяются на обучение *с учителем* и *без учителя*. Методы обучения с учителем учитывают зависимости между входными переменными и откликами, то есть плотности условного распределения $p(y | \mathbf{x})$. Цель методов обучения без учителя — аппроксимировать структуру или закономерности во входных данных, то есть смоделировать плотность безусловного распределения $p(\mathbf{x})$. В обучении без учителя не делается никаких предположений о переменных отклика, управляющих обучением модели, и для поиска зависимо-

стей между признаками или точками данных анализируется только матрица плана. Каноническим примером обучения без учителя является задача кластеризации, которую можно определить как задачу группировки точек данных в кластеры так, чтобы точки данных в одном кластере были похожи друг на друга, а в разных — отличались. Эта задача не требует переменных отклика, но опирается на меру сходства точек данных. Обучение без учителя широко используется в маркетинговых приложениях для исследования и анализа данных. Кластеризация профилей клиентов и интерпретация полученных результатов, например, один из важнейших методов маркетинговой аналитики. Однако в программных приложениях мы больше занимаемся автоматизацией, чем исследованием и интерактивным анализом. Обучение представлением — это одно из применений методов обучения без учителя, которые могут пригодиться в этом контексте, поэтому здесь сосредоточимся на аспектах обучения представлению, а не на обучении без учителя в целом.

2.5.1. Метод главных компонент

Метод главных компонент (Principal Component Analysis, PCA) — мощный прием поиска сжатого и некоррелированного представления данных. Метод главных компонент — это математический метод, который преобразует данные в новое представление, обладающее определенными свойствами, а также создает артефакты, описывающие структуру данных. В зависимости от приложения нас могут интересовать разные свойства преобразования методом главных компонент, поэтому в следующих разделах мы последовательно обсудим их, но имейте в виду, что все они основаны на одном алгоритме.

2.5.1.1. Устранение корреляции

В маркетинговых приложениях данные обычно соответствуют наблюдаемым входным данным, свойствам и выходным данным некоторых реальных маркетинговых процессов. Примерами таких процессов могут служить маркетинговые кампании, отношения между клиентами и продуктами, а также связь цены и спроса. Каждый признак в матрице плана можно рассматривать как сигнал, несущий информацию о процессе. Мы не обладаем полным знанием процесса и наблюдаем только его определенные *проекции* на размерности признаков, доступные во входных данных, подобно тому как физический объект можно сфотографировать с разных точек. Например, мы не наблюдаем вкусы и мысли потребителей напрямую, а фиксируем определенные сигналы, например покупки, которые частично отражают вкусы, мысли и решения. Представления, полученные таким способом, вероятно, будут иметь некоторую избыточность, а размерности, вероятно, будут коррелированы, так же как фотографии одного и того же объекта с разных точек являются избыточными и коррелированными. Эта идея проиллюстрирована на рис. 2.7.

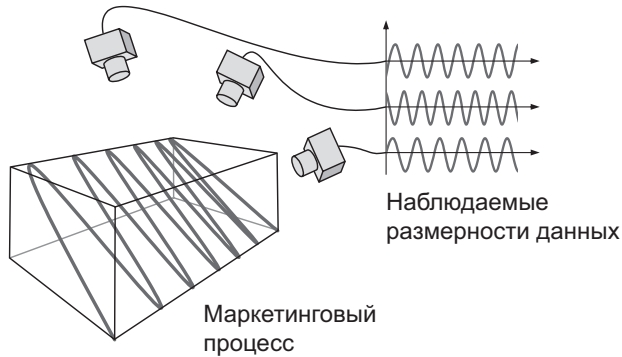


Рис. 2.7. Различные размерности данных могут коррелироваться, потому что являются проекциями одного и того же реального процесса

Возникает задача поиска нового, потенциально меньшего набора статистически независимых признаков, обеспечивающих менее избыточное и более структурированное представление данных. Для ее решения можно применить метод главных компонент, предположив, что статистическую независимость признаков можно заменить нулевой корреляцией, что может быть верно или неверно, в зависимости от распределения данных. Это ограничительное предположение, потому что для статистической независимости с нулевой корреляцией значения признаков должны быть совместно нормально распределены. Если распределение отличается, метод главных компонент может не достигнуть цели устранения корреляции (но все еще оставаться полезным из-за других свойств). Рассмотрим матрицу \mathbf{X} с размерами $n \times m$, где n — число точек данных, а m — число признаков. Предположим, что данные центрированы, то есть $E[\mathbf{x}] = 0$. Если это не так, можно вычесть среднее из всех точек данных. Тогда задачу декорреляции можно определить в терминах ковариационной матрицы \mathbf{C} :

$$\mathbf{C}_x = \text{Var}[\mathbf{x}] = \frac{1}{n-1} \mathbf{X}^T \mathbf{X}. \quad (2.63)$$

Это квадратная симметричная матрица $m \times m$, в которой диагональные элементы являются дисперсиями соответствующих признаков, а недиагональные — ковариациями. Признаки считаются некоррелированными, если ковариационная матрица диагональна, то есть если все недиагональные элементы в \mathbf{C} равны нулю. Если ковариационная матрица не диагональна, значит, признаки коррелированы, что затрудняет понимание распределения \mathbf{x} , потому что его нельзя описать в терминах распределений отдельных признаков. Наш следующий шаг — найти преобразование \mathbf{X} , которое создаст другую матрицу плана \mathbf{Z} с размерами $n \times m$, для которой ковариационная матрица диагональна. В методе главных компонент

предполагается, что такую матрицу можно получить с помощью линейного преобразования:

$$Z = XT, \quad (2.64)$$

где T — матрица преобразования. Чтобы получить эту матрицу, покажем сначала, как можно разложить матрицу плана на векторы, соответствующие направлениям максимальной дисперсии в матрице плана.

Сначала найдем направления максимальной дисперсии в данных. Их можно рассматривать как основные оси облака точек на диаграмме рассеяния, как показано на рис. 2.8. Каждое из направлений можно определить как вектор, поэтому начнем с нахождения m -мерного единичного вектора, удовлетворяющего следующему условию:

$$v_1 = \operatorname{argmax}_v \|Xv\|^2, \quad \|v\| = 1. \quad (2.65)$$

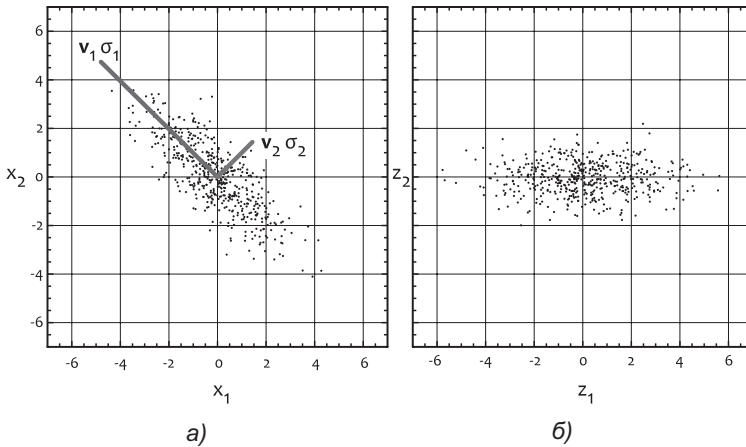


Рис. 2.8. Пример применения метода главных компонент. а) Набор данных из 500 нормально распределенных точек и соответствующих главных компонент. Признаки x_1 и x_2 сильно коррелируют. б) Декоррелированное представление, полученное методом главных компонент. Признаки z_1 и z_2 не коррелированы

Этот вектор соответствует оси с максимальной дисперсией. Затем найдем второй единичный вектор, ортогональный первому, объясняющий оставшуюся дисперсию:

$$v_2 = \operatorname{argmax}_v \|Xv\|, \quad \|v\| = 1 \text{ и } v \cdot v_1 = 0. \quad (2.66)$$

Продолжаем этот процесс, требуя, чтобы каждый следующий вектор был ортогонален всем предыдущим, и, предполагая, что матрица плана имеет ранг r , мы можем создать r ненулевых векторов \mathbf{v} . В силу своих свойств каждый вектор соответствует оси с максимальной оставшейся дисперсией в матрице плана. Эти векторы называются *главными компонентами* матрицы плана.

Обозначим матрицу $m \times r$, составленную из векторов-столбцов \mathbf{v} как \mathbf{V} . Поскольку все созданные нами единичные векторы ортогональны, эта матрица является ортонормированной по столбцам, то есть

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}. \quad (2.67)$$

Единичные векторы \mathbf{v} фиксируют направления дисперсии, но не величину. Посчитаем эти значения отдельно и обозначим их как

$$\sigma_i = \|\mathbf{X} \mathbf{v}_i\|. \quad (2.68)$$

Каждая главная компонента объясняет только *оставшуюся* дисперсию, поэтому первая компонента имеет наибольшее значение, а значения всех последующих уменьшаются:

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r. \quad (2.69)$$

Обозначим диагональную матрицу $r \times r$ со значениями σ на главной диагонали как

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r). \quad (2.70)$$

На данный момент у нас имеется ортонормированный базис векторов \mathbf{V} и соответствующие коэффициенты масштабирования Σ . Чтобы завершить разложение матрицы плана, нам нужен третий фактор, который проецирует матрицу на базис главных компонент или, говоря иначе, подмешивает базис в матрицу плана. Обозначив этот фактор как \mathbf{U} , мы можем определить разложение как

$$\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T, \quad (2.71)$$

откуда можно получить \mathbf{U} , решив уравнение 2.71 относительно этого фактора:

$$\mathbf{U} = \mathbf{X} \mathbf{V} \Sigma^{-1}. \quad (2.72)$$

Разложение, определяемое уравнением 2.71, известно как сингулярное разложение (Singular Value Decomposition, SVD). Подводя итог, главные компоненты, то есть

столбцы матрицы V , можно интерпретировать как ортогональные главные оси, соответствующие направлениям дисперсии. Столбцы $U\Sigma$ можно интерпретировать как коэффициенты, которые в сочетании с главными векторами производят X . Также можно показать, что матрица U является ортонормированной по столбцам:

$$U^T U = I. \quad (2.73)$$

Получив сингулярное разложение матрицы плана, мы можем использовать факторы для поиска линейного преобразования T , устраняющего корреляцию. Рассмотрим произведение

$$Z = XV \quad (2.74)$$

и вычислим его ковариационную матрицу, используя тот факт, что матрицы V и U ортогональны:

$$\begin{aligned} C_z &= \frac{1}{n-1} Z^T Z = \\ &= \frac{1}{n-1} V^T X^T X V = \\ &= \frac{1}{n-1} V^T V \Sigma^2 V^T V = \\ &= \frac{1}{n-2} \Sigma^2. \end{aligned} \quad (2.75)$$

Так как Σ^2 является диагональной матрицей, представление Z не коррелировано. Это означает, что декоррелирующее преобразование, которое мы искали, фактически задается матрицей V . Это преобразование по сути является матрицей поворота, потому что матрица главных компонент ортонормирована. Отметим, что предположение о линейности преобразования является довольно ограничительным. В примере, приведенном на рис. 2.8, это дает хороший результат, потому что набор данных имеет эллиптическую форму, что является результатом нормального распределения; в этом случае корреляцию между признаками можно устранить простым поворотом. С наборами данных с более сложной формой ситуация может сложиться совсем иначе.

2.5.1.2. Уменьшение размерности

Ключевым свойством метода главных компонент является упорядоченность главных векторов по величине объясняемой дисперсии. Это свойство очень важно, потому что позволяет утверждать, что размерности в данных с высокой дисперсией, как правило, более информативны и несут более сильный сигнал. Например, можно сказать, что оси x_1 и x_2 на рис. 2.8 одинаково важны, и одномерные представления,

полученные при проецировании точек данных на ось x_1 или x_2 , являются плохими приближениями исходного двумерного представления. В то же время размерности z_1 и z_2 в некоррелированном наборе данных, полученные с помощью метода главных компонент, не одинаково важны и размерность z_2 можно отбросить, пожертвовав относительно небольшой долей информации.

Этому свойству можно найти разные применения. Первое такое применение — уменьшение размерности данных, например, когда m -мерный набор данных \mathbf{X} необходимо свести к k -мерному и $k < m$. В методе главных компонент такую свертку можно выполнить с помощью усеченной матрицы \mathbf{V}_k , включающей только первые k векторов главных компонент, вычислив новое представление данных как

$$\mathbf{Z}_k = \mathbf{X}\mathbf{V}_k. \quad (2.76)$$

Матрица \mathbf{Z}_k имеет только k столбцов, соответствующих первым главным осям. Этот подход часто используется в визуализации данных для отображения многомерного набора данных в двух- или трехмерное пространство, которое можно показать на графике.

Второе важное применение — аппроксимация матрицы плана *матрицей меньшего ранга*. Рассмотрим сингулярное разложение, заданное выражением 2.71. По аналогии с уменьшением размерности из матрицы \mathbf{V} можно удалить наименее значимые столбцы и создать усеченную версию \mathbf{V}_k . Факторы \mathbf{U} и Σ также нужно усечь, чтобы удалить столбцы и строки, соответствующие наименее значимым компонентам, как показано на рис. 2.9. Обозначим эти усеченные версии как \mathbf{U}_k и Σ_k соответственно. После этого матрицу плана \mathbf{X} можно аппроксимировать произведением этих усеченных факторов:

$$\hat{\mathbf{X}}_k = \mathbf{U}_k \Sigma_k \mathbf{V}_k. \quad (2.77)$$

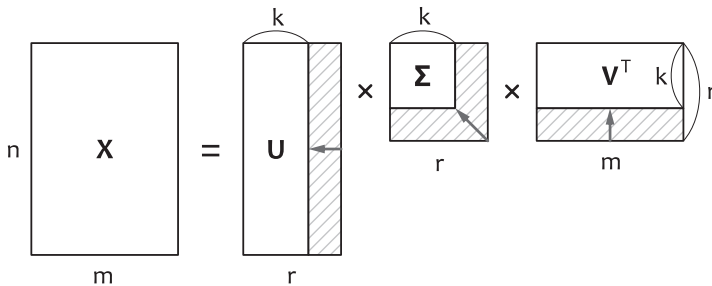


Рис. 2.9. Снижение размерности удалением осей, объясняющих меньшие доли дисперсии

Восстановленная матрица \hat{X} имеет тот же размер, что и X , но, конечно же, страдает некоторой погрешностью аппроксимации из-за отброшенных размерностей. Эта погрешность невелика, если отбросить только наименее значимое измерение, и увеличивается по мере удаления большего числа столбцов из факторов сингулярного разложения. Фактически можно показать, что такой метод аппроксимации матрицы является оптимальным в том смысле, что реконструированная матрица имеет наименьшую возможную погрешность аппроксимации с учетом ограничения, требующего, чтобы ранги факторов не превышали k . Иными словами, выражение 2.77 является решением следующей задачи оптимизации:

$$\min_A \|X - A\| \quad (2.78)$$

с условием $\text{rank}(A) \leq k$.

Аппроксимация матрицей меньшего ранга с успехом используется в маркетинговых приложениях, особенно в услугах поиска и выбора рекомендаций, поскольку позволяет справиться с разреженными, зашумленными и избыточными наборами данных. Часто к ним относятся наборы данных, описывающие связи между двумя сущностями. Например:

- Матрица плана может отражать связь между клиентами и продуктами. Каждая строка матрицы соответствует клиенту, каждый столбец представляет продукт, а каждый элемент — это метрика связи, например количество покупок. На практике такая матрица почти наверняка будет очень разреженной, поскольку каждый пользователь покупает только небольшую часть доступных продуктов. Кроме того, данные в наборе будут сильно коррелированы, потому что многие продукты похожи друг на друга, так же как многие клиенты имеют схожие покупательские привычки.
- В приложениях поиска такие тексты, как описания продуктов, часто моделируются в виде векторов, каждый элемент в которых соответствует слову; как следствие, длина вектора равна общему числу разных слов в словаре. Коллекцию таких текстов можно представить в виде матрицы, где каждая строка соответствует текстовому документу, а каждый столбец — слову. Эта матрица может быть разреженной, особенно для коротких текстов, и избыточной, потому что слова с определенным семантическим значением часто появляются вместе.

В примерах выше каждый элемент x_{ij} матрицы плана является некоторой мерой сходства между сущностями, например сходства между клиентом и продуктом или между словом и документом. Однако исходные значения сходства в матрице плана часто являются зашумленными и неполными. Можно ли создать более гладкую модель, предсказывающую сходство для любой пары сущностей? Один

из возможных подходов представить сходство как скалярное произведение двух векторов

$$\hat{x}_{ij} = \mathbf{p}_i \cdot \mathbf{q}_j^T, \quad (2.79)$$

где первая сущность (например, клиент) представлена некоторым числовым вектором \mathbf{p} , а вторая (например, продукта) — числовым вектором \mathbf{q} . Длины векторов k обычно выбираются значительно меньше размера матрицы плана. Переписав это выражение в матричной форме

$$\hat{\mathbf{X}} = \mathbf{P} \cdot \mathbf{Q}^T, \quad (2.80)$$

можно заметить, что векторы \mathbf{p} и \mathbf{q} , минимизирующие среднюю ошибку аппроксимации сходства, на самом деле можно получить из выражения аппроксимации матрицей меньшего ранга 2.77:

$$\begin{aligned} \mathbf{P} &= \mathbf{U}_k \Sigma_k \\ \mathbf{Q} &= \mathbf{V}_k. \end{aligned} \quad (2.81)$$

Это очень важный результат, потому что помогает преобразовать разреженные и избыточные представления сущностей в компактные и плотные числовые векторы. Это очень мощный и универсальный метод моделирования, который широко используется в следующих главах.

2.5.2. Кластеризация

Кластеризация — это процесс группировки похожих элементов. Ее также можно рассматривать как деление набора данных на кластеры, внутри которых точки данных имеют большое сходство, а точки в разных кластерах — малое.

В традиционных подходах к маркетингу кластеризация чаще используется как метод исследовательского анализа данных. Посредством кластеризации набор данных можно разделить на относительно небольшое число кластеров, и каждый кластер можно описать, интерпретировать и изучить как один объект. Каноническим примером является *сегментация* клиентов, когда большое количество профилей клиентов делится на несколько кластеров, или сегментов, на основе меры сходства, учитывающей демографические особенности, поведение и покупательские привычки. Каждый такой сегмент можно описать на основе его геометрического центра в векторном пространстве профилей (среднего профиля) и разброса значений признаков, то есть типичное описание кластера может выглядеть, например, так: *экономные клиенты в возрасте до 30 лет, которых можно привлечь в основном*

через цифровые каналы. То есть сегментация позволяет обобщить большой набор данных в несколько точек, пригодных для анализа вручную. Тема сегментирования является одной из крупнейших и наиболее стратегических тем в маркетинговой аналитике, потому что большую часть корпоративной маркетинговой стратегии можно построить на сегментах клиентов, их типичных потребностях и свойствах.

Программные приложения, более ориентированные на выполнение, чем на стратегический анализ, часто используют результаты сегментации в качестве дополнительных признаков. Например, вектор профиля клиента, содержащий такие значения, как возраст, доход и среднемесячные расходы, можно дополнить метками сегментов, описывающими человека как *любителя выгодных покупок, модника, приверженного определенному бренду*, и т. д. Эти дополнительные признаки, как и любые другие, можно использовать в предиктивном моделировании для повышения точности прогнозирования и интерпретируемости результатов. С этой точки зрения кластеризацию можно рассматривать как метод проектирования признаков.

Кластеризацию также можно применять к сущностям, которые используются в качестве признаков при моделировании других сущностей. В онлайн-рекламе, например, профили пользователей обычно включают URL-адреса, посещавшиеся пользователем в прошлом, поэтому данные профиля могут выглядеть так:

$$\begin{aligned}\text{user 1: } & (url_1, url_2, url_3, \dots), \\ \text{user 2: } & (url_4, url_5, url_6, \dots).\end{aligned}$$

Такое представление может быть очень разреженным из-за очень большого числа различных URL-адресов. Следовательно, легко может случиться так, что профили пользователей будут иметь очень мало общих URL-адресов или вообще не иметь их, поэтому модель, использующая такое представление профиля как признак, не сможет точно соответствовать данным. Впрочем, каждый URL-адрес можно связать с вектором атрибутов, таких как доменное имя и родственные веб-сайты, и применить алгоритм кластеризации для объединения URL-адресов в категории. Категории, созданные алгоритмом кластеризации, могут иметь некоторое семантическое значение. Например, один кластер URL-адресов может соответствовать преимущественно спортивным ресурсам, тогда как другой — технологическим. В результате появляется возможность сопоставить URL-адреса в профилях пользователей с кластерами и выразить профили с точки зрения поведенческих особенностей:

$$\begin{aligned}\text{user 1: } & (\text{спорт, мода, технологии, } \dots), \\ \text{user 2: } & (\text{новости, мода, спорт, } \dots).\end{aligned}$$

В предположении, что число кластеров намного меньше числа URL-адресов, такое представление гораздо плотнее, и многие профили будут иметь общие метки. Такое использование кластеризации является примером обучения представлением, направленного на поиск удобного набора признаков, и существенно отличается от исследовательского анализа.

Кластеризация — сложная задача, потому что обучение происходит без учителя, и, следовательно, цель оптимизации нельзя определить однозначно. Существует несколько семейств алгоритмов кластеризации, реализующих разные подходы к решению этой задачи. Один из таких подходов основан на интерпретации кластеризации как задачи обучения модели: поскольку цель состоит в группировке похожих элементов, можно предположить, что каждая группа является результатом некоторого случайного и неизвестного процесса, и можно подобрать такую *смесь* распределений, которая достаточно точно (в смысле максимального правдоподобия) аппроксимирует наблюдаемый набор данных. Этот подход проиллюстрирован на рис. 2.10, где кластеры определяются путем подбора смеси трех нормальных распределений данных.

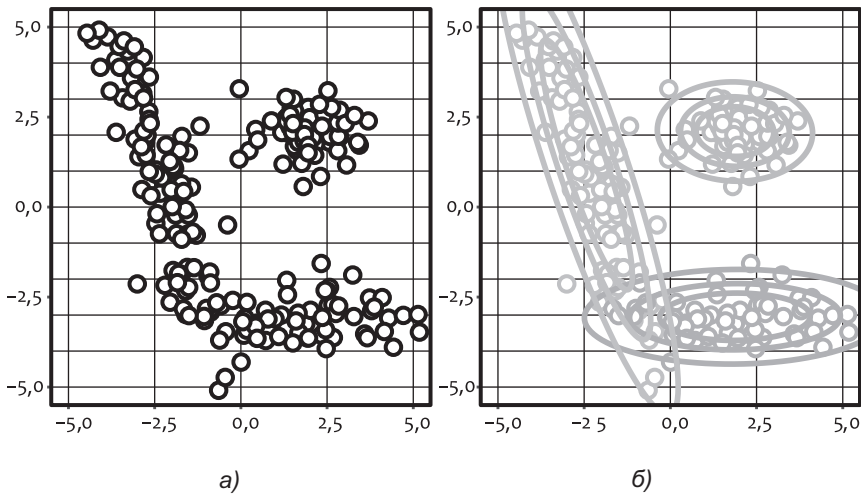


Рис. 2.10. Пример кластеризации с использованием моделирования смеси распределений. а) Исходный набор данных. б) Три кластера, найденных подбором смеси из трех нормальных распределений

Эту модель можно задать так:

$$p(\mathbf{x}) = \sum_{k=1}^K w_k N(\mathbf{x} | \mu_k, \Sigma_k), \quad (2.82)$$

где K — число кластеров, w_k — веса смешивания, а μ_k и Σ_k — матрица средних и ковариационная матрица распределений в смеси соответственно (обратите внимание, что нормальное распределение — это только один из возможных вариантов, и можно использовать смесь других распределений). Имея параметры распределений, каждую точку данных легко можно отнести к соответствующему кластеру, опираясь на плотности вероятностей в этой точке. Верно и обратное — зная кластеры, к которым относятся точки данных, легко можно оценить параметры каждого распределения в смеси. Проблема, однако, в том, что нам не известно ни то, как точки отображаются в кластеры, ни параметры распределений; известны только исходные точки данных. Это влечет усложнение функций правдоподобия, которые гораздо сложнее вычислить, чем вероятности, изучаемые в предыдущих разделах. Однако есть ряд итерационных методов, способных найти приближенное решение. Наиболее широко используются алгоритмы максимизации ожиданий (Expectation-Maximization EM) и K -средних.

2.6. Более специализированные модели

Стандартные методы обучения с учителем и без учителя отвечают наиболее типичным потребностям моделирования в маркетинге. Многие маркетинговые задачи относительно просто преобразовать в такие стандартные задачи предиктивного моделирования. Однако некоторые маркетинговые задачи могут потребовать применения более специализированных методов анализа или сложных экономических моделей, объединяющих бизнес-цель и основные примитивы предиктивного моделирования. Некоторые из этих методов первоначально были разработаны в экономике, другие — в теории игр, биологии и социальных науках. В этом разделе мы познакомимся с несколькими специализированными моделями и методами, расширяющими набор стандартных методов машинного обучения.

2.6.1. Теория потребительского выбора

Понимание и прогнозирование потребительского выбора — одна из самых фундаментальных проблем в маркетинге, а также в экономике в целом, поскольку многие важные вопросы, связанные с проектированием изделий, планированием ассортимента и распространением, нельзя решить, если не изучить спрос достаточно хорошо. В этом разделе мы рассмотрим проблему дискретного выбора, то есть ситуацию, когда принимающий решение сталкивается с выбором из множества альтернатив. Например, потребитель решает, какой из конкурирующих товаров купить, отменить ли подписку на определенную услугу или нет и т. д. Можно предположить, что принимающий решение сравнивает варианты по-

следовательно и непротиворечиво (если вариант k является предпочтительнее варианта m , а вариант m предпочтительнее n , то k предпочтительнее n), и выбирает наиболее предпочтительный, тогда допустимо ввести виртуальную числовую меру, значение которой пропорционально *ценности* данного варианта для принимающего решение.

Обозначим ценность для принимающего решение n в случае выбора варианта j из набора альтернатив $(1, \dots, J)$ как Y_{nj} . Принимающий решение выбирает вариант Y_{nj} среди других вариантов, если $Y_{nj} > Y_{ni}$ для всех случаев $i \neq j$. Ценность одного и того же варианта j не обязательно одинакова для всех принимающих решение из-за различий во вкусах, доходах и другие личных особенностей.

Модель потребительского выбора можно создать на основе известных свойств индивидов и альтернатив. Однако каждый принимающий решение почти наверняка будет учитывать какие-то дополнительные особенности, влияющие на выбор, которые не известны создателю модели. Говоря более формально, можно утверждать, что ценность Y_{nj} является функцией известных \mathbf{x}_{nj} и скрытых \mathbf{h}_{nj} факторов:

$$Y_{nj} = Y(\mathbf{x}_{nj}, \mathbf{h}_{nj}). \quad (2.83)$$

Скрытые факторы \mathbf{h}_{nj} известны принимающему решение, но не создателю модели, поэтому модель ценности $V_{nj} = V(\mathbf{x}_{nj})$ аппроксимирует истинную ценность Y_{nj} с некоторой ошибкой ε_{nj} , которую можно рассматривать как случайную величину:

$$Y_{nj} = V_{nj} + \varepsilon_{nj}. \quad (2.84)$$

Такой подход к анализу ценности известен как *случайная модель ценности* (random utility model). Определение 2.84 позволяет выразить вероятность выбора альтернативы j принимающим решение n как

$$\begin{aligned} P_{ni} &= \Pr(Y_{ni} > Y_{nj}, \forall j \neq i) = \\ &= \Pr(V_{ni} + \varepsilon_{ni} > V_{nj} + \varepsilon_{nj}, \forall j \neq i). \end{aligned} \quad (2.85)$$

Обозначим случайный вектор ошибок как

$$\varepsilon_n = (\varepsilon_{n1}, \dots, \varepsilon_{nj}). \quad (2.86)$$

Предположив, что распределение ε_n известно, можно оценить вероятность выбора интегрированием плотности вероятности $p(\varepsilon_n)$:

$$P_{nj} = \int_{\varepsilon} \mathbb{I}(V_{ni} + \varepsilon_{ni} > V_{nj} + \varepsilon_{nj}, \forall j \neq i) p(\varepsilon_n) d\varepsilon_n, \quad (2.87)$$

где \mathbb{I} — индикаторная функция, равная 1 для истинного аргумента, и 0 — для ложного. Нам не требуется делать каких-то конкретных предположений относительно модели V , чтобы вычислить выражение 2.87, — мы свободны в выборе любой линейной или нелинейной функции известных факторов \mathbf{x}_{nj} оценки ценности. Однако, исходя из соображений практичности, мы должны сделать некоторые предположения относительно распределения $p(\epsilon_n)$, того, чтобы сделать оценку P_{nj} пригодной для анализа.

Разные предположения об остаточных ошибках ϵ_n приводят к разным моделям выбора с разными достоинствами и ограничениями. Конечная цель состоит в том, чтобы найти более разумную с вычислительной точки зрения формулу, выражающую P_{nj} как функцию V_{nj} , которая в свою очередь является функцией наблюдаемых свойств \mathbf{x}_{nj} и некоторых параметров $\boldsymbol{\omega}$. Это может быть, например, линейная модель:

$$V_{nj} = \boldsymbol{\omega}^T \mathbf{x}_{nj}. \quad (2.88)$$

Поскольку P_{nj} обычно можно оценить по известным статистикам для альтернатив, которые были доступны в прошлом, также можно оценить параметры $\boldsymbol{\omega}$ и построить предиктивную модель для P_{nj} . Эту модель можно использовать для оценки новых альтернатив, определенных в терминах свойств \mathbf{x} , и, следовательно, предсказывать экономические показатели, такие как спрос на новый продукт. В следующем разделе мы обсудим одну из самых простых, но мощных и практичных моделей — *полиномиальную модель вероятности с логистическим распределением* (полиномиальная logit-модель), — чтобы показать, как можно получить для P_{nj} выражения с разумной вычислительной сложностью. Эта модель будет использоваться в последующих главах как компонент более сложных моделей для конкретных маркетинговых задач.

2.6.1.1. Полиномиальная модель с логистическим распределением

Полиномиальную logit-модель (Multinomial Logit Model, MNL) можно получить из случайной модели ценности, предположив, что остаточные ошибки ϵ_{nj} независимы и подчиняются распределению Гумбеля. Предположение о распределении Гумбеля приводит к удобной модели и может также рассматриваться как практическая аппроксимация нормального распределения [Train, 2003]. Предположение о независимых распределениях является гораздо более строгим и приводит к ограничениям, которые будут рассматриваться далее в этом разделе. В общем случае распределение Гумбеля используется для описания распределения максимального или минимального значений в ряду случайных точек данных, взятых из некоторого базового распределения. Например, если сгенерировать пакеты случайных

чисел, извлекая их из нормального распределения, а затем из каждого пакета взять максимальное значение, распределение этих максимумов можно смоделировать с помощью распределения Гумбеля. Этот подход с успехом применяется для моделирования некоторых экстремальных событий, таких как землетрясения, дефекты производства и отказы оборудования. Например, производитель, выпускающий медикаменты партиями, может использовать распределение Гумбеля для моделирования вероятности производства партии, где уровень активных химических компонентов выше максимально допустимого уровня. Плотность вероятности распределения Гумбеля определяется как

$$p(\varepsilon_{nj}) = e^{-\varepsilon_{nj}} \exp(-e^{-\varepsilon_{nj}}). \quad (2.89)$$

Кумулятивное распределение задается как

$$F(\varepsilon_{nj}) = \exp(-e^{-\varepsilon_{nj}}). \quad (2.90)$$

Чтобы воспользоваться предположением о том, что остаточные ошибки подчиняются распределению Гумбеля, сначала перепишем вероятность выбора, заданную уравнением 2.85, следующим образом:

$$\begin{aligned} P_{ni} &= \Pr(V_{ni} + \varepsilon_{ni} > V_{nj} + \varepsilon_{nj}, \forall j \neq i) = \\ &= \Pr(\varepsilon_{nj} < \varepsilon_{ni} + V_{ni} - V_{nj}, \forall j \neq i). \end{aligned} \quad (2.91)$$

Предположив на мгновение, что ε_{ni} дано и используя независимость ошибок, мы можем заявить, что

$$P_{ni} | \varepsilon_{ni} = \prod_{j \neq i} \Pr(\varepsilon_{nj} < \varepsilon_{ni} + V_{ni} - V_{nj}). \quad (2.92)$$

Члены справа фактически являются кумулятивными распределениями ε_{nj} , поэтому, вставив определение распределения Гумбеля 2.90, мы получим:

$$P_{ni} | \varepsilon_{ni} = \prod_{j \neq i} \exp(-e^{-(\varepsilon_{ni} + V_{ni} - V_{nj})}). \quad (2.93)$$

Поскольку значение ε_{ni} фактически не задано, интегрируем его плотность вероятности, чтобы получить полное выражение для P_{ni} :

$$P_{ni} = \int_{\varepsilon} (P_{ni} | \varepsilon_{ni}) \cdot e^{-\varepsilon_{ni}} \exp(-e^{-\varepsilon_{ni}}) d\varepsilon_{ni}. \quad (2.94)$$

Краткое выражение в замкнутой форме для вероятности выбора можно получить непосредственно из уравнения выше, применив алгебраические преобразования,

которые мы опустим здесь ради краткости; результатом является каноническая формула для модели MNL:

$$P_{ni} = \frac{e^{V_{ni}}}{\sum_j e^{V_{nj}}}. \quad (2.95)$$

Модель, представленная уравнением 2.95, имеет несколько важных свойств и следствий.

НЕЗАВИСИМОСТЬ ОТ ПОСТОРОННИХ АЛЬТЕРНАТИВ. Один из самых важных вопросов, на которые необходимо ответить при моделировании выбора, — как ценность одной из альтернатив влияет на другие альтернативы. Например, производитель может быть заинтересован в оценке доли потребителей, которых можно отвлечь от конкурентов снижением цены на продукцию или запуском производства нового продукта. Модель MNL подразумевает, что увеличение или уменьшение вероятности одной альтернативы равномерно влияет на все другие альтернативы. Чтобы увидеть это, рассмотрим соотношение любых двух вероятностей:

$$\frac{P_{ni}}{P_{nj}} = \frac{e^{V_{ni}} / \sum_k e^{V_{nk}}}{e^{V_{nj}} / \sum_k e^{V_{nk}}} = \frac{e^{V_{ni}}}{e^{V_{nj}}} = e^{V_{ni} - V_{nj}}. \quad (2.96)$$

Отношение вероятностей зависит только от отношения соответствующих ценностей — свойство, обычно называемое *независимостью от посторонних альтернатив*, — при изменении ценности V_{ni} попарные отношения для всех других пар P_{np}/P_{nq} остаются неизменными. Это свойство MNL является несколько ограничивающим, потому что продукты внутри исследуемой группы не всегда являются полностью взаимозаменяемыми и могут иметь место более сложные шаблоны замещения. Это ограничение можно проиллюстрировать парадоксом Debreu, 1960.

Рассмотрим транспортную систему, в которой потребитель выбирает между автомобилем и автобусом и где вероятности первоначального выбора равны:

$$P_{car} = P_{bus} = 1/2. \quad (2.97)$$

Предположим теперь, что появился второй автобус, идентичный первому. Допустим, что разница только в цвете: первый красный, а второй синий. Модель MNL равномерно перераспределит вероятности

$$P_{car} = P_{red\ bus} = P_{blue\ bus} = 1/3, \quad (2.98)$$

поскольку оба автобуса имеют одинаковую ценность. Однако реалистичнее было бы предположить, что соотношение

$$P_{car} / (P_{red\ bus} + P_{blue\ bus} + \dots) \quad (2.99)$$

останется постоянным, независимо от количества одинаковых автобусов, что дает вероятности $P_{car} = 1/2$ и $P_{red\ bus} = P_{blue\ bus} = 1/4$.

ПОЛНОТА МОДЕЛИ ЦЕННОСТИ. Независимость от остаточных ошибок ε_{ni} подразумевает, что модель ценности V_{ni} должна охватывать все факторы, влияющие на выбор. Если модель V_{ni} не полна, некоторые систематические смещения начинают просачиваться в компоненты ошибок и нарушают предположение о независимости. Например, мы можем построить модель ценности для стиральной машины, используя цену p и потребление энергии s как предиктивные переменные, то есть $V_{ni} = w_1 p_i + w_2 s_i$. Однако выбор, скорее всего, будет зависеть от дохода потребителя g , который может быть неизвестен, поэтому ценность фактически будет иметь вид

$$Y_{ni} = w_1 p_i + w_2 s_i + w_3 g_n + \varepsilon_{ni} = V_{ni} + \varepsilon_{ni}^*, \quad (2.100)$$

где $\varepsilon_{ni}^* = w_3 g_n + \varepsilon_{ni}$ и представляет ошибки, не являющиеся независимыми из-за случайной переменной g_n .

ПРЕДЕЛЬНАЯ ВЕРОЯТНОСТЬ ВЫБОРА. Вероятность выбора является сигмоидной функцией ценности V_{ni} , как показано на рис. 2.11. Это означает, что небольшие изменения ценности вызывают существенное увеличение или уменьшение вероятности, что соответствующий вариант будет выбран, только если вероятность близка к 0,5, то есть если принимающий решение находится в предельном состоянии. Если вероятность выбора той или иной альтернативы мала или велика, даже крупные изменения ценности оказывают ограниченное влияние на вероятность.

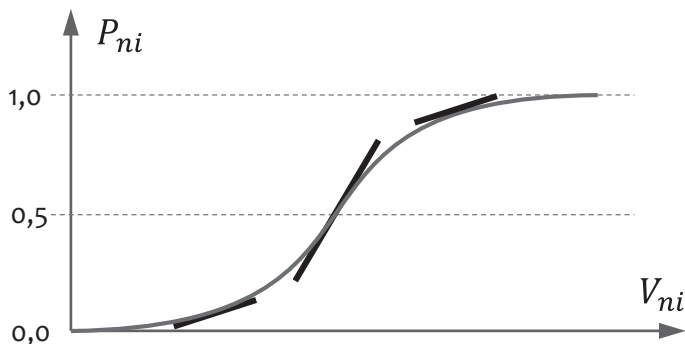


Рис. 2.11. Вероятность выбора как функция ценности и ее производные в разных точках

Такая S-образная зависимость между V_{ni} и P_{ni} предполагает, что инвестиции должны быть сосредоточены на разработке альтернатив с промежуточной вероятностью выбора большинством лиц, принимающих решения. Например, онлайн-ритейлер, совершенствующий свою сеть доставки заказов, может рассчитывать на самую высокую отдачу от инвестиций, улучшая услуги в областях, где он уже имеет среднюю долю рынка, в то время как области с очень низкими или очень высокими долями, вероятно, будут менее восприимчивы к улучшениям.

2.6.1.2. Оценка полиномиальной logit-модели

Посмотрим теперь, как можно оценить параметры модели ценности V_{ni} по обучающему набору данных. Предположим, что нам известны признаки \mathbf{x}_{ni} , включенные в модель ценности для некоторого подмножества лиц, принимающих решения $n = 1, \dots, N$ и каждого варианта $i = 1, \dots, J$, а также фактически сделанный выбор. Пусть $y_{ni} \in \{0, 1\}$ — наблюдаемый выбор лица n , принявшего решение относительно варианта i ; он равен 1, если принимающий решение выбрал эту альтернативу, и 0 в противном случае. Согласно предположению о независимости остаточных ошибок, мы можем выразить вероятность фактического выбора как

$$\prod_i (P_{ni})^{y_{ni}}. \quad (2.101)$$

Вероятность, что все, принимающие решения, в данном наборе данных сделали выбор, который мы фактически наблюдаем, то есть правдоподобие набора данных можно выразить как

$$L(\mathbf{w}) = \prod_{n=1}^N \prod_i (P_{ni})^{y_{ni}} \quad (2.102)$$

при условии, что все решения независимы. Следовательно, логарифмическое правдоподобие будет иметь вид

$$\begin{aligned} LL(\mathbf{w}) &= \sum_{n=1}^N \sum_i y_{ni} \log(P_{ni}) = \\ &= \sum_{n=1}^N \sum_i y_{ni} \log \frac{e^{V_{ni}}}{\sum_k e^{V_{nk}}}, \end{aligned} \quad (2.103)$$

где V_{ni} является функцией от \mathbf{w} и \mathbf{x} , например линейная модель $V_{ni} = \mathbf{w}^T \mathbf{x}_{ni}$. Логарифмическое правдоподобие, описываемое уравнением 2.103, можно далее вычислить, взяв градиент относительно \mathbf{w} и применив численные методы оптимизации.

2.6.2. Анализ выживаемости

Методы классификации, даже самые простые, такие как логистическая регрессия, — это мощный инструмент для оценки вероятностей действий потребителей. Например, вероятность отклика на рекламное электронное письмо можно оценить, построив модель, использующую в качестве признаков атрибуты клиента, такие как количество покупок, и в качестве метки ответа — признак ответа клиента на предыдущее рекламное письмо. Этот подход широко используется на практике, но имеет несколько недостатков, которые мы подробно рассмотрим в последующих главах. Во-первых, во многих маркетинговых приложениях удобнее и эффективнее оценивать время до события, а не вероятность события. Например, для маркетинговой системы может быть полезнее оценить время до следующей покупки или время до отказа от подписки, чем вероятность этих событий. Во-вторых, маркетинговые данные очень часто включают записи с неизвестными или пропущенными результатами, которые нельзя должным образом учесть в моделях классификации. Возвращаясь к примеру с отказом от подписки, зачастую невозможно отличить клиентов, которые не сбежали, от клиентов, которые *пока* не сбежали, потому что мы строим предиктивную модель на определенный момент времени и не можем ждать бесконечно, пока появятся окончательные результаты для всех клиентов. Следовательно, мы знаем только результаты для клиентов, которые сбежали, и только их уверенно можем отметить как отрицательные образцы; остальные записи являются неполными, и нет никакой уверенности, что эти клиенты не сбегут в будущем, поэтому можно утверждать, что маркировка их как положительные или отрицательные образцы на самом деле не является действительной. Из этого следует, что было бы неправильно использовать модель классификации с бинарной переменной на выходе, определяемой на основе результатов, наблюдаемых в настоящий момент, и нам нужна другая статистическая структура для решения такого рода задачи.

Для медицинских и биологических исследований в свое время была разработана комплексная платформа для моделирования времени-до-события и обработки неполных данных. Основное внимание в исследованиях уделялось выживаемости людей после врачебной помощи, поэтому данную платформу назвали анализом выживаемости. Давайте опишем основные методы этой платформы, начав с базовой терминологии. Главная цель анализа выживаемости — оценить время до интересующего события и количественно объяснить, как это время зависит от параметров лечения, индивидуальных особенностей пациентов и других независимых переменных. В маркетинге аналогом лечения можно считать стимулирование или побуждение, например посредством рекламы. Роль события обычно играет покупка, активация, отказ от подписки или любое другое действие клиента, на которое маркетолог хотел бы повлиять. Отметим, что положительным результатом

лечения может быть приближение или отдаление момента наступления события, в зависимости от конечной цели. Рекламные объявления, например, направлены на стимулирование более ранних покупок, тогда как поощрительные предложения направлены на отдаление события отказа от подписки. В медицинских исследованиях, напротив, время обычно измеряется от диагноза до смерти, поэтому стандартная терминология анализа выживаемости предполагает, что событие соответствует некоторому отрицательному результату, что может сбивать с толку, когда имеет место обратное.

Как отмечалось выше, некоторые события могут быть неизвестны, в том смысле что результаты не наблюдались во время исследования. Эти пробелы в результатах могут иметь место из-за того, что результат не был известен на момент анализа (клиент еще не сделал покупку, но может сделать ее в будущем) или запись с информацией о клиенте была потеряна (например, из-за истечения срока действия cookie в веб-браузере). Записи с неизвестными результатами называют *цензурированными*. К моменту проведения анализа мы изначально имеем набор наблюдений, каждое из которых имеет время стимулирования и, не обязательно, время события.

Время между стимулированием и событием называют *временем выживаемости*. Мы можем преобразовать исходные наблюдения, упорядочив по времени стимулирования, чтобы наблюдения для k физических лиц (клиентов) представить в виде последовательности пар:

$$(t_1, \delta_1), \dots, (t_k, \delta_k), \quad t_1 \leq \dots \leq t_k, \quad (2.104)$$

где t обозначает временную метку события, а δ — индикатор, равный 1, если наблюдение не цензурировано, и 0 в противном случае. Обычно предполагается непрерывность временной шкалы, но у двух клиентов может быть одинаковое время события, поэтому входные данные можно обобщить как

$$(t_1, d_1), \dots, (t_n, d_n), \quad (2.105)$$

где n — число различных времен событий; d_i — общее число событий, наблюдаемых в момент времени t_i . Мы также предполагаем, что события не повторяются, то есть для каждого человека может произойти не более одного события. Для многих маркетинговых событий, таких как покупки, это предположение не является верным в буквальном смысле, но обычно мы можем обойти его, создавая отдельные модели для первого, второго и последующих событий, как мы обсудим в последующих главах. На данный момент нас интересует только распределение событий, и мы не пытаемся объяснить зависимость между временем выживаемости и параметрами стимулирования или особенностями клиентов.

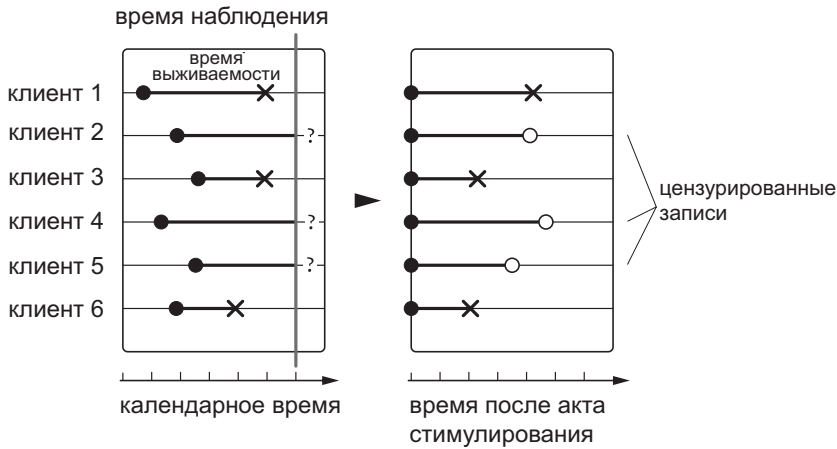


Рис. 2.12. Подготовка для анализа выживаемости. Закрашенные кружки соответствуют актам стимулирования. Крестики — соответствующим событиям. Незакрашенные кружки обозначают цензурированные записи

2.6.2.1. Функция выживаемости

Распределение времени выживаемости можно описать в терминах вероятности выживаемости, *функции выживаемости* $S(t)$, которая определяется как вероятность, что человек проживет от начального момента до момента времени t . Функция выживаемости — фундаментальная характеристика, описывающая динамику группы клиентов. Если функция выживаемости падает резко, событие с большой долей вероятности относительно быстро наступит для большинства клиентов. Если функция падает медленно, почти наверняка для большинства клиентов событие наступит в относительно отдаленной точке в будущем.

Обозначим время выживаемости клиента как T и его функцию плотности вероятности как $f(t)$. Кумулятивная функция распределения времени выживаемости, соответствующая вероятности события до момента времени t , будет тогда равна

$$F(t) = \Pr(T \leq t) = \int_0^t f(\tau) d\tau, \quad (2.106)$$

а функцию выживаемости можно определить как

$$S(t) = \Pr(T > t) = 1 - F(t). \quad (2.107)$$

Значение функции выживаемости в момент времени t соответствует доле клиентов, для которых событие еще не наступило. Обратите внимание, что статистические

свойства времени выживаемости, такие как среднее, медиана и доверительные интервалы, можно оценить на основе кумулятивной функции распределения. Следовательно, эти свойства можно получить при наличии оценки функции выживаемости.

Функцию выживаемости можно вычислить по наблюдаемым данным с учетом как цензурированных, так и не цензурированных записей, исходя из предположения о том, что события независимы друг от друга. В этом случае кумулятивную вероятность выживаемости можно получить путем умножения вероятности выживаемости из одного интервала на вероятность в следующем. Более формально вероятность дожить до времени t можно оценить прямолинейно:

$$S_t = \frac{n_t - d_t}{n_t} = 1 - \frac{d_t}{n_t}, \quad (2.108)$$

где n_t — количество людей, для которых событие еще не наступило в момент времени t , а d_t — число людей, для которых событие уже наступило к моменту t . Путем перемножения вероятностей от начального момента времени до момента t можно оценить суммарную вероятность выжить, то есть функцию выживаемости:

$$\hat{S}(t) = \prod_{i \leq t} \left(1 - \frac{d_i}{n_i} \right). \quad (2.109)$$

Эта оценка известна как оценка Каплана–Мейера (Kaplan–Meier), и можно доказать, что она является оценкой максимального правдоподобия [Kaplan and Meier, 1958]. Функция выживаемости равна 1 в нулевой момент времени, а затем, с увеличением времени, каждая точка данных вносит свой вклад в оценку. Проиллюстрируем оценку функции выживаемости на небольшом числовом примере.

ПРИМЕР 2.1

Предположим, что мы анализируем группу из 14 клиентов после того, как каждый из них получил рекламное письмо. Все письма были отправлены в разное время, и в ходе анализа фиксировалось время до первой покупки, прошедшее после отправки письма. Набор наблюдаемых данных выглядит следующим образом:

$$\begin{aligned} t &= \{2, 3, 3, 3, 4, 6, 7, 8, 12, 12, 14, 15, 20, 23\}, \\ \delta &= \{1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1\}, \end{aligned} \quad (2.110)$$

где i -й элемент в множестве t — наблюдаемое время события для i -го клиента, измеренное в днях с момента отправки электронного письма. Множество δ содержит признак цензурированности наблюдения (0) или нецензурированности (1). Например, первый клиент совершил покупку на второй день после получения письма, а третий не совершил покупки к моменту анализа, хотя получил письмо за три дня до даты анализа. В этом контексте под вероятностью выживаемости подразумевается вероятность не совершить покупку к данному моменту времени. Многократно применяя формулу 2.109, получаем следующую последовательность:

$$\begin{aligned} S(0) &= 1 \quad (\text{в начальный момент времени все клиенты «живы»}) \\ S(2) &= 1 - \frac{1}{14} = 0,93 \\ S(3) &= S(2) \cdot \left(1 - \frac{2}{13}\right) = 0,79 \\ &\dots \end{aligned} \tag{2.111}$$

Этот результат соответствует ступенчатой *кривой выживаемости*, изображенной на рис. 2.13. Кривая выживаемости обобщает динамику группы клиентов, а кроме того, допустимо сравнивать кривые для различных групп. Например, можно построить кривые выживаемости для клиентов, которых стимулировали в рамках рекламной акции, и для тех, кому рекламные письма не рассылались, и графически оценить эффективность рекламной акции.

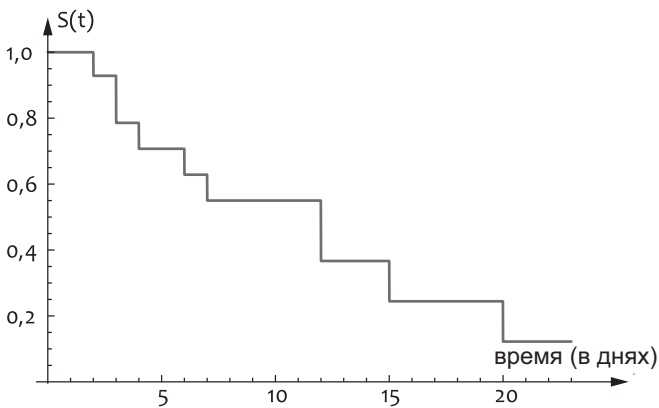


Рис. 2.13. Оценка функции выживаемости для набора данных в определении 2.110

2.6.2.2. Функция риска

Второе важное понятие в анализе выживаемости — *функция риска*. Если функция выживаемости фокусируется на вероятности, что событие не произойдет, то есть на *выживании*, то функция риска описывает риск наступления события. Как мы увидим позже, эта перспектива удобна для анализа влияния различных факторов, такие как параметры лечения, на время выживаемости.

Функция риска $h(t)$ определяется как мгновенная скорость риска, то есть вероятность события в бесконечно малом промежутке времени между t и $t+dt$ с учетом того, что человек дожил до времени t :

$$h(t) = \lim_{dt \rightarrow 0} \frac{\Pr(t < T \leq t + dt \mid T > t)}{dt}. \quad (2.112)$$

Функция риска имеет определенную связь с функцией выживаемости. Чтобы увидеть это, давайте сначала разложим условную вероятность в определении 2.112 на два фактора; отметьте, что один из них соответствует функции выживаемости:

$$\begin{aligned} h(t) &= \lim_{dt \rightarrow 0} \frac{\Pr(t < T \leq t + dt)}{dt \cdot \Pr(T > t)} = \\ &= \lim_{dt \rightarrow 0} \frac{\Pr(t < T \leq t + dt)}{dt \cdot S(t)} = \\ &= \lim_{dt \rightarrow 0} \frac{F(t + dt) - F(t)}{dt \cdot S(t)}. \end{aligned} \quad (2.113)$$

Затем вспомним, что функция плотности вероятности определяется как

$$f(t) = \lim_{dt \rightarrow 0} \frac{F(t + dt) - F(t)}{dt}. \quad (2.114)$$

Подставив это определение, а также определение функции выживаемости 2.107, мы получим следующий результат:

$$\begin{aligned} h(t) &= \frac{f(t)}{S(t)} = \frac{f(t)}{1 - F(t)} = \\ &= -\frac{d}{dt} \log(1 - F(t)) = \\ &= -\frac{d}{dt} \log(S(t)). \end{aligned} \quad (2.115)$$

Решая это уравнение относительно $S(t)$, можно выразить функцию выживаемости как функцию $h(t)$:

$$S(t) = \exp(-H(t)), \quad (2.116)$$

где

$$H(t) = \int_0^t h(\tau) d\tau = -\log(S(t)) \quad (2.117)$$

называется кумулятивной функцией риска. Эта прямая связь позволяет переключаться между функциями риска и выживаемости в анализе.

2.6.2.3. Регрессионный анализ выживаемости

Базовые функции выживаемости и риска можно использовать для описания поведения группы клиентов или сравнения различных групп друг с другом. Но этого недостаточно для случаев, когда требуется понять и предсказать, как на выживаемость и риск влияют такие факторы, как маркетинговые действия и индивидуальные особенности клиента. Эта задача аналогична задачам классификации и регрессии, в том смысле что время выживаемости должно предсказываться как функция наблюдаемых факторов, то есть независимых переменных.

Предположим, что каждому индивидууму соответствует вектор \mathbf{x} , состоящий из p независимых переменных. То есть каждый индивидуум представлен тремя значениями:

t — время выживаемости, или цензурированное время;

δ — признак цензурирования, принимает значение 1 для наблюдаемых событий и 0 — для цензурированных;

\mathbf{x} — вектор признаков.

Набор входных данных содержит наблюдения для k клиентов:

$$(t_1, \delta_1, \mathbf{x}_1), \dots, (t_k, \delta_k, \mathbf{x}_k), \quad t_1 \leq \dots \leq t_k. \quad (2.118)$$

В маркетинговых приложениях вектор признаков может включать демографические и поведенческие свойства клиента, маркетинговые коммуникации с ним и т. д. Цель состоит в том, чтобы определить и обучить модель, которая выражает функции выживаемости и риска как функцию от \mathbf{x} . Поскольку $S(t)$ и $h(t)$ являются вероятностями, мы можем построить разные регрессионные модели выживаемо-

сти, предполагая разные распределения вероятностей и разные функциональные зависимости между признаками \mathbf{x} и параметрами распределения.

Чаще всего в роли регрессионных моделей выживаемости используются модели *пропорциональных рисков*. Это семейство моделей основано на предположении, что единичное увеличение наблюдаемых факторов мультипликативно по отношению к степени риска, то есть

$$h(t | \mathbf{w}, \mathbf{x}) = h_0(t) \cdot r(\mathbf{w}, \mathbf{x}), \quad (2.119)$$

где $h_0(t)$ является базовым риском, r — отношение рисков, увеличивающее или уменьшающее базовый риск в зависимости от факторов, а \mathbf{w} — вектор параметров модели. Обратите внимание, что базовый риск зависит не от конкретного человека, а от отношения рисков. Другими словами, отношение рисков определяет, как свойства индивидуума, закодированные в векторе признаков, влияют на уровень риска. Отношение рисков не может быть отрицательным, так как коэффициент риска не отрицателен, поэтому он обычно моделируется как экспоненциальная функция:

$$h(t | \mathbf{w}, \mathbf{x}) = h_0(t) \cdot \exp(\mathbf{w}^T \mathbf{x}). \quad (2.120)$$

Эту модель можно считать линейной по отношению к логарифму отношения рисков индивидуума к базовому риску:

$$\log r(\mathbf{w}, \mathbf{x}) = \log \frac{h(t | \mathbf{x})}{h_0(t)} = \mathbf{w}^T \mathbf{x}. \quad (2.121)$$

Что касается базового риска $h_0(t)$, у нас на выбор есть два варианта: параметрический и непараметрический. Параметрический подход предполагает, что риск подчиняется определенному распределению. В этом случае мы получаем полностью параметрическую модель, которую необходимо обучить на исходных данных путем поиска оптимальных значений параметров \mathbf{w} и параметров распределения. Недостатком этого подхода является предположение, что базовый риск изменяется во времени определенным образом, поэтому мы должны быть уверены, что выбранное распределение соответствует данным. С другой стороны, непараметрический подход сглаживает зашумленные данные и обеспечивает простую модель базового риска.

Второй вариант заключается в использовании непараметрической модели базового риска, которую можно получить на основе данных с использованием оценки Каплана–Мейера или других методов. Это приводит к полупараметрической модели для общего риска, где параметрическая часть определяется выражением 2.120, а базовый риск $h_0(t)$ представляет непараметрическую часть. Это решение известно

как модель пропорциональных рисков Кокса [Cox, 1972]. Преимущество модели Кокса, как мы увидим ниже, заключается в возможности оценить коэффициенты риска без необходимости вычислять базовую функцию риска или делать какие-либо предположения о структуре базового риска. Это делает ее очень удобной для применений, где требуется определить только факторы риска, а не абсолютные значения. Недостатком модели Кокса является необходимость определения базового риска с помощью параметрических методов. Важно также иметь в виду, что модель Кокса относится к семейству моделей пропорциональных рисков и, следовательно, основана на предположении о пропорциональности рисков, что может быть верно или неверно для наблюдаемых данных. Модель Кокса широко используется во многих областях, включая маркетинг, и мы будем использовать ее в качестве основного инструмента для анализа выживаемости в следующей главе.

Наш следующий шаг — определить параметры модели Кокса по данным. Стандартное решение этой задачи состоит в том, чтобы получить правдоподобие модели, а затем найти параметры, максимизирующие его. Проблема, однако, в том, что наблюдения могут оказаться цензурированными, а это требует от нас указать, как такие записи должны учитываться в определении правдоподобия. Прежде всего отметим, что каждое наблюдение вносит свой вклад в величину подобия. Если i -е наблюдение цензурировано, это означает вероятность дожить до t_i :

$$L_i(\boldsymbol{w}) = S(t_i | \boldsymbol{w}, \boldsymbol{x}). \quad (2.122)$$

Если наблюдение не цензурировано, оно вносит вклад в вероятность возникновения события в момент t_i , которая определяется с помощью функции плотности вероятности времени выживаемости:

$$L_i(\boldsymbol{w}) = f(t_i) = h(t_i | \boldsymbol{w}, \boldsymbol{x}) S(t_i | \boldsymbol{w}, \boldsymbol{x}). \quad (2.123)$$

То есть полное правдоподобие можно выразить так:

$$L_i(\boldsymbol{w}) = \prod_{i=1}^k h(t_i | \boldsymbol{w}, \boldsymbol{x})^{\delta_i} S(t_i | \boldsymbol{w}, \boldsymbol{x}). \quad (2.124)$$

Мы не сможем максимизировать это выражение с использованием численных методов, не указав форму базового риска. Однако можно аппроксимировать полное правдоподобие с помощью другой меры, называемой частичным правдоподобием. Для начала введем понятие *группы риска* в момент времени t , которое определяется как множество лиц с риском наступления события в момент t , то есть лиц, для которых событие еще не наступило:

$$R(t) = \{i : t_i \geq t\}. \quad (2.125)$$

В целях упрощения предположим также отсутствие связей между событиями, то есть что все времена событий t_i различны¹. В этом случае *частичное правдоподобие* можно определить с помощью условной вероятности, что конкретный индивидум i из группы риска на момент t_i потерпит неудачу к этому моменту и произойдет ровно одна неудача [Сох, 1972, 1975]. Эта вероятность задается областью под кривой риска для небольшого временного интервала dt , поэтому правдоподобие, добавленное индивидуумом i , можно выразить как

$$L_i(\mathbf{w}) = \frac{h(t_i | \mathbf{w}, \mathbf{x}) dt}{\sum_{j \in R(t_i)} h(t_j | \mathbf{w}, \mathbf{x}_j) dt}. \quad (2.126)$$

Подставив определение модели Кокса 2.120 в это частичное правдоподобие, можем увидеть, что базовые риски отменяют друг друга, и мы получаем:

$$L_i(\mathbf{w}) = \frac{\exp(\mathbf{w}^T \mathbf{x}_i)}{\sum_{j \in R(t_i)} \exp(\mathbf{w}^T \mathbf{x}_j)}. \quad (2.127)$$

Наконец, частичное правдоподобие для всего набора обучающих данных является произведением отдельных частичных вероятностей, заданных уравнением 2.127:

$$L(\mathbf{w}) = \prod_{i=1}^k \left[\frac{\exp(\mathbf{w}^T \mathbf{x}_i)}{\sum_{j \in R(t_i)} \exp(\mathbf{w}^T \mathbf{x}_j)} \right]^{\delta_i}. \quad (2.128)$$

Эта модель правдоподобия не зависит от функции риска, поэтому ее можно обучить с помощью численных методов относительно весов \mathbf{w} . Это, в свою очередь, позволяет получить отношения рисков, определяемые уравнением 2.121. Возможность получить отношение рисков без вычисления функции риска является одним из ключевых преимуществ модели Кокса.

До сих пор все наше внимание было сосредоточено на получении коэффициентов регрессии. Теперь мы должны определить, как вычислить функции базового риска и выживаемости. Прежде всего отметим, что ожидаемое число событий в момент времени t_i можно аппроксимировать площадью под функцией риска на небольшом интервале между t_i и $t_i + dt$:

¹ Случай со связанными событиями сложнее, однако есть ряд обобщений, помогающих учесть связи [Breslow, 1974; Efron, 1977]. В большинстве маркетинговых приложений от связей можно избавиться, распределяя противоречивые наблюдения с небольшим отрывом.

$$\hat{d}_i = \sum_{j \in R(t_i)} h_0(t_i) \exp(\mathbf{w}^T \mathbf{x}_j) dt. \quad (2.129)$$

Это отношение можно переписать иначе:

$$\hat{h}_0(t_i) dt = \frac{\hat{d}_i}{\sum_{j \in R(t_i)} \exp(\mathbf{w}^T \mathbf{x}_j)}, \quad (2.130)$$

а затем аппроксимировать кумулятивную функцию риска:

$$\hat{H}_0(t) = \sum_{i < t} \hat{h}_0(t_i) dt. \quad (2.131)$$

Этот результат известен как оценка Бреслоу [Breslow, 1972]. Она позволяет оценить базовую функцию выживаемости подстановкой оценки в выражение 2.116:

$$\hat{S}_0(t) = \exp(-\hat{H}_0(t)). \quad (2.132)$$

Наконец, полную функцию выживаемости можно получить непосредственно из определения модели Кокса в уравнении 2.120 и выражения 2.116:

$$\begin{aligned} S(t | \mathbf{x}) &= \exp \left[- \int_0^t h_0(\tau) \exp(\mathbf{w}^T \mathbf{x}) d\tau \right] = \\ &= \exp \left[- \int_0^t h_0(\tau) d\tau \right]^{\exp(\mathbf{w}^T \mathbf{x})} = \\ &= S_0(t)^{\exp(\mathbf{w}^T \mathbf{x})}. \end{aligned} \quad (2.133)$$

Функции выживаемости для разных значений признаков \mathbf{x} можно построить вместе и сравнить друг с другом, чтобы оценить влияние разных признаков на распределение времени выживаемости. В следующей главе мы продолжим изучение этого и других практических применений анализа выживаемости в рекламе и маркетинговых коммуникациях.

2.6.3. Теория аукционов

Как говорилось в главе 1, алгоритмический подход способствует развитию маркетинговых услуг, которые могут быть предложены биржами. Биржа или брокер любого другого типа между поставщиком услуг и клиентом добавляет дополнительный уровень сложности, потому что кроме достижения основных маркетинговых целей оба, поставщик и клиент, должны оптимизировать свои стратегии покупки и продажи услуг.

Основная цель биржи услуг заключается в поддержании конкуренции между покупателями за ограниченный ресурс, такой как места для размещения рекламы. Стандартный способ решения этой задачи — организация *аукциона*, где каждый покупатель делает *ставку*, и ресурс, выставленный на аукцион, достается участнику с максимальной ставкой. Однако правила аукциона можно настроить по-разному, поэтому мы потратим некоторое время на изучение типов аукционов.

Во-первых, важно понимать, что потенциальные покупатели участвуют в аукционе, потому что для каждого из них аукционный ресурс имеет определенную стоимость, и они стремятся получить прибыль, стараясь приобрести желаемое по как можно более низкой цене. Следовательно, для участника торгов критически важно правильно оценить стоимость ресурса, и, как результат, мы можем классифицировать все аукционы по следующим типам стоимости.

ЧАСТНАЯ СТОИМОСТЬ. Каждый участник оценивает ресурс независимо от других, и его оценка не зависит от других предложений, даже если они известны.

ВЗАИМОЗАВИСИМАЯ СТОИМОСТЬ. Фактическая стоимость неизвестна участникам, и, хотя каждый имеет собственную оценку, информация о других предложениях может помочь улучшить ее. Например, участник торгов, высоко оценивший ресурс, может уменьшить ставку, если другие уменьшат свои ставки, поскольку эта дополнительная информация может указывать на наличие негативных факторов, не известных данному покупателю, но каким-то образом ставших известными другим.

ОБЩАЯ СТОИМОСТЬ. Это частный случай аукциона с взаимозависимой стоимостью, когда фактическая стоимость одинакова для всех участников. В качестве примеров аукционов с общей стоимостью можно назвать продажу природных ресурсов, таких как нефть или древесина; финансовых активов, таких как облигации; компаний. Во всех этих случаях истинная стоимость может быть не известна точно во время аукциона, и участники должны оценить ее на основе ограниченной информации, которой располагают, но в конечном итоге стоимость станет известна (фактическое количество извлекаемой нефти, долгосрочные результаты деятельности компании и т. д.) и она одинакова для всех участников.

Хотя стоимость часто в той или иной степени взаимозависима, способность участника торгов воспользоваться преимуществами знания других ставок зависит от процесса аукциона. Ниже перечислены четыре основных типа аукционов, изученных в теории и применяемых на практике:

ОТКРЫТЫЕ ТОРГИ. Участники торгов имеют информацию о величине ставок, сделанных другими.

- Открытый повышающий (английский) аукцион. Торги начинаются с низкой цены, и затем она постепенно повышается. Участник может в любой момент покинуть аукцион. Когда остается только один участник, аукцион завершается и победивший платит окончательную цену.
- Открытый понижающий (голландский) аукцион. Торги начинаются с высокой цены, и затем она постепенно понижается. Аукцион завершается, когда обнаруживается участник, готовый заплатить текущую цену.

ЗАКРЫТЫЕ ТОРГИ. Участники торгов не имеют информации о величине ставок, сделанных другими.

- Закрытый аукцион первой цены. Все участники одновременно делают ставки так, что ни один не знает ставок других участников. Выигрывает участник, давший самую высокую цену, и платит ее.
- Закрытый аукцион второй цены (аукцион Викри). Так же как на аукционе первой цены, все участники одновременно делают ставки, и выигрывает участник, давший самую высокую цену, но платит по второй по величине ставке.

Задача оптимизации для открытых аукционов может показаться динамичной, но в действительности она статична и эквивалентна задаче оптимизации закрытых аукционов. Голландский аукцион заканчивается после первой ставки, поэтому участники не получают никакой дополнительной информации в процессе торгов и могут принять решение о цене заранее. То есть голландский аукцион эквивалентен закрытому аукциону первой цены в том смысле, что независимо от стратегии, выбранной участником, в нем используется та же исходная информация и победитель определяется, как в аукционах с частной и взаимозависимой стоимостью. В английском аукционе с частной стоимостью участник также может оценить товар заранее. На каждом этапе торгов и роста цены участник должен сравнить текущую ставку со своей оценкой и либо сделать новую ставку, полученную из текущей прибавлением некоторой величины, либо выйти из аукциона, если цена превысила его оценку. Следовательно, английский аукцион эквивалентен закрытому аукциону второй цены для частной стоимости, хотя это не относится к взаимозависимым стоимостям, потому что в случае английского аукциона участник может извлечь уроки из наблюдаемых ставок.

Теперь подробнее остановимся на аукционе Викри, чтобы получить инструменты для построения моделей оптимизации, которые включают аукционы. Мы сосредоточимся на аукционе Викри, потому что он удобнее для анализа и широко используется в практических приложениях, впрочем, похожие результаты можно получить для других типов аукционов, используя более продвинутые методы анализа.

Во-первых, докажем, что оптимальной стратегией для участников торгов является истинная стоимость. Взгляните на рис. 2.14, где участник торгов оценивает товар по цене v , но делает более низкую ставку $v - \delta$. Если вторая по величине ставка другого участника равна p , возможны следующие три результата:

1. $p > v$: участник проигрывает; в этом случае не имеет значения, какую ставку он сделал, v или $v - \delta$.
2. $p < v - \delta$: участник выигрывает и платит цену p ; в этом случае не имеет значения, какую ставку он сделал, v или $v - \delta$.
3. $v - \delta < p < v$: участник проигрывает; ставка v означала бы выигрыш и маржу $v - p$.

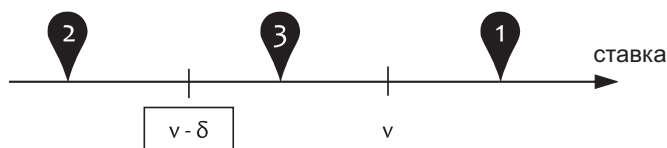


Рис. 2.14. Аукцион Викри — ставки ниже истинной стоимости

Таким образом, ставка ниже истинной стоимости всегда дает тот же или худший результат, что и ставка с истинной стоимостью. На рис. 2.15 показана противоположная ситуация, когда ставка выше истинной стоимости. В этом случае у нас снова три возможных исхода:

1. $p > v + \delta$: участник проигрывает; в этом случае не имеет значения, какую ставку он сделал, v или $v + \delta$.
2. $p < v$: участник выигрывает и платит цену p ; в этом случае не имеет значения, какую ставку он сделал, v или $v + \delta$.
3. $v < p < v + \delta$: участник выигрывает и платит цену p , потеряв $p - v$, тогда как ставка v означала бы проигрыш без финансовых потерь.

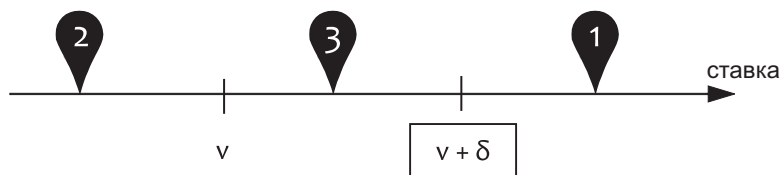


Рис. 2.15. Аукцион Викри — ставки выше истинной стоимости

Можно сделать вывод, что оптимальной стратегией являются торги по истинной стоимости. Этот простой результат предполагает, что при рассмотрении параметров маркетинга для закрытых аукционов мы должны сосредоточиться на оценке ожидаемых доходов участника.

Теперь встанем на точку зрения продавца аукциона и оценим его доход. Предположим, что в аукционе участвуют n участников, предложивших цены V_1, \dots, V_n — независимые и одинаково распределенные случайные величины из некоторого распределения $F(v)$ с плотностью вероятности $f(v)$. Помня, что статистика k -го порядка $V(k)$ образца равна k -му наименьшему значению, ожидаемый доход (*revenue*) можно выразить как среднее значение статистики второго сверху порядка, которое соответствует второй наивысшей ставке:

$$\text{revenue} = \mathbb{E}[V_{n-1}]. \quad (2.134)$$

Рассмотрим срез функции плотности вероятности для порядковых статистик:

$$\Pr(v < V_{(k)} < v + dv). \quad (2.135)$$

Это вероятность, что $k - 1$ ставок в выборке из n окажутся меньше v , ровно одна заявка попадет в диапазон $[v, v + dv]$, а оставшиеся $n - k$ ставок будут выше v . Эти три условия можно выразить с помощью функций кумулятивного распределения $F(v)$ и плотности вероятности $f(v)$, откуда получаем следующее выражение плотности вероятности для порядковых статистик:

$$\begin{aligned} f(V_{(k)}) &= \lim_{dv \rightarrow 0} \Pr(v < V_{(k)} < v + dv) = \\ &= \binom{n}{k-1} [F(v)]^{k-1} \cdot (n-k+1)f(v) \cdot [1-F(v)]^{n-k} = \\ &= \frac{n!}{(k-1)!(n-k)!} f(v) [F(v)]^{k-1} [1-F(v)]^{n-k}. \end{aligned} \quad (2.136)$$

Мы можем упростить это выражение, сделав некоторые предположения относительно распределения ставок. Например, если ставки равномерно распределены между нулем и единицей, выражение сокращается до

$$f(V_{(k)}) = \frac{n!}{(k-1)!(n-k)!} v^{k-1} (1-v)^{n-k}. \quad (2.137)$$

Это бета-распределение, поэтому для получения ожидаемого дохода аукциониста можно использовать стандартную формулу:

$$\mathbb{E}[V_{(n-1)}] = \frac{n-1}{n+1}. \quad (2.138)$$

Этот результат согласуется с интуитивным ожиданием, что увеличение числа участников торгов приведет к общему увеличению доходов. Мы воспользуемся этими результатами в следующих главах, в основном для оптимизации процесса торгов на биржах. Однако стоит отметить, что другие маркетинговые процессы, в том числе такие важные, как отбор потребителем лучших предложений на рынке, также можно смоделировать как аукционы, что делает теорию аукционов важным инструментом для создания программных решений.

2.7. Итоги

- Многие маркетинговые задачи можно выразить в виде задач оптимизации, в которых предметом оптимизации является бизнес-результат, а переменными — бизнес-действия.
- Зависимость между действиями и бизнес-результатами часто можно получить из исторических данных. Эту задачу можно решить с помощью методов обучения с учителем.
- Основной целью обучения с учителем является оценка условного распределения отклика по имеющимся исходным данным. Во многих практических приложениях эту задачу можно свести к поиску наиболее вероятных исходов. Двумя основными типами задач обучения с учителем являются классификация и регрессия.
- Число параметров предиктивной модели может быть фиксированным или увеличиваться с размером обучающего набора данных. Модели первого типа называют параметрическими, а второго — непараметрическими.
- Обучение модели можно рассматривать как задачу оптимизации, где требуется подобрать параметры модели, максимизирующие вероятность следования наблюдаемых данных распределению модели.
- Многие задачи обучения с учителем можно решить с помощью линейных моделей, то есть моделей, которые определяют либо зависимость между входом и выходом в виде линейной функции, либо границу между классами в виде прямой линии. Наиболее простыми примерами линейных моделей являются линейная и логистическая регрессия.

- Нелинейные зависимости и границы решения можно аппроксимировать с помощью нелинейных моделей. Примерами методов нелинейного моделирования могут служить ядерные методы, деревья решений и нейронные сети.
- Маркетинговые данные могут иметь избыточную структуру, поскольку разные функции и метрики являются проекциями одного и того же маркетингового процесса. Структура может быть неоптимальной для анализа и моделирования, а лучшее представление данных может быть найдено путем удаления корреляций, уменьшения размерности данных и кластеризации точек данных и сущностей. Некоторые из этих задач можно решить с помощью методов обучения без учителя, таких как анализ главных компонент и кластеризация.
- Некоторые маркетинговые задачи не решаются с помощью стандартных методов машинного обучения и требуют применения более специализированных моделей и методов. Примерами таких моделей могут служить модели потребительского выбора, методы анализа выживаемости и теория аукционов.

3

Продвижение и реклама

Каждый товар или услуга имеют свой целевой рынок — группу потребителей, на которую они нацелены. Различие между целевыми и нецелевыми группами часто нечеткое и неоднозначное, так как потребители различаются по доходам, покупательскому поведению, лояльности к бренду и многим другим параметрам. Разнообразие клиентов зачастую настолько велико, что предложение, созданное для среднестатистического потребителя, то есть для каждого, ничьим потребностям не соответствует. Поэтому для бизнеса критически важно идентифицировать наиболее релевантных потребителей и адаптировать свои предложения на основе особенностей этих потребителей. Эта задача возникает практически во всех маркетинговых приложениях и особенно важную роль играет в рекламе и продвижении, потому что эффективность этих услуг напрямую зависит от умения правильно определить аудиторию и послать правильный сигнал.

Задачу поиска оптимального соответствия между потребителями и предложениями часто можно рассматривать с двух точек зрения. Во-первых, ее можно обозначить как задачу поиска правильных предложений для конкретного клиента. Это задача *поиска продукта*, и мы обсудим ее в следующих главах, посвященных поиску и рекомендациям. С другой точки зрения задача состоит в поиске подходящих клиентов для данного предложения. Эта задача известна как *таргетирование* (выбор целевой группы) и является основной темой данной главы. Имейте в виду, что грань между поиском продукта и таргетированием услуг проводится, главным образом, по видам применения (интерактивный просмотр или прямая реклама), а методологии, используемые для реализации услуг, иногда можно рассматривать с обеих точек зрения. Рассмотрим в качестве примера сегментацию клиентской базы. Можно утверждать, что сегментация сначала определяет правильные группы клиентов, а затем под каждый сегмент производится адаптация предложений. Однако сегментацию можно также рас-

смаатривать как метод распределения различных предложений среди наиболее подходящих клиентов.

Несмотря на то что таргетирование связано с сопоставлением клиентов и предложений, его не следует рассматривать просто как набор методов для проведения соединительных линий между этими двумя сущностями. Напротив, его следует рассматривать как задачу оптимизации качества обслуживания клиентов, которая обусловлена сочетанием нескольких бизнес-целей и контролирует множество разных маркетинговых мероприятий. Цель программной системы состоит в том, чтобы развернуть эти начальные цели в подробный план выполнения и определенные правила, которые можно использовать для управления взаимодействиями с клиентами.

Мы начнем эту главу с обзора среды продвижения розничной торговли, которая поможет нам лучше понять задачу таргетирования. Затем опишем основу таргетирования продвижения, включая более формальное определение бизнес-целей, основные строительные блоки поведенческого моделирования и более сложные конструкции, используемые в маркетинговых кампаниях. Затем мы обсудим среду онлайн-рекламы и связанные с ней методы таргетирования. Даже при том что среда розничной торговли и онлайн-рекламы дополняют друг друга и многие методы таргетирования универсальны, мы будем изучать их отдельно из-за серьезных структурных различий в целях. Наконец, мы обсудим способы оценки эффективности методов таргетирования и маркетинговых кампаний. Оценка играет чрезвычайно важную роль во всех маркетинговых приложениях, и основа, которую мы разработаем, также будет использоваться в других программных услугах, включая поиск, рекомендации и ценообразование. В этой главе мы постараемся не затрагивать оптимизацию цен, хотя это и важная часть продвижения. Эта тема будет рассмотрена далее в отдельной главе.

3.1. Среда

Первая бизнес-среда, которую мы рассмотрим, — это продвижение продаж, которое широко используется в розничной торговле и управлении отношения к бренду. Цель продвижения состоит в увеличении ценности продукта для потребителей, чтобы улучшить продажи или построить лучшие отношения. Рекламные акции с целью продвижения могут проводиться производителями продуктов, поставщиками услуг или ретейлерами. В некоторых рыночных вертикалях, таких как *производство товаров широкого потребления*, производители и ретейлеры часто сотрудничают в рекламных кампаниях, при этом производитель покрывает прямые затраты на кампанию, а ретейлер предоставляет физические и цифровые каналы для передачи предложений аудитории. Как говорилось в главе 1, такие взаимо-

действия между промоутерами и владельцами потребительской базы могут быть благоприятными для алгоритмического подхода, поэтому мы будем использовать это обстоятельство в качестве основы для дальнейшего изучения. Однако большинство методов, которые мы будем разрабатывать, не ограничиваются средой производства товаров широкого потребления и могут применяться в других областях, таких как телекоммуникации или страхование.

Модель среды стимулирования продаж представлена на рис. 3.1. Основные элементы этой среды, допущения и терминологию можно описать следующим образом.

- *Потребитель* — это любой человек, потребляющий товары. *Клиент* — это человек, приобретающий что-то у фирмы. Наконец, *перспективный клиент* — это человек, еще не являющийся клиентом, но уже известный фирме в том смысле, что фирма может общаться с ним (например, если человек зарегистрировался на сайте и оставил свой электронный адрес). Потребителями мы также называем всех, кто взаимодействует с онлайн-каналами в качестве пользователей.
- *Производители (или бренды)* производят продукты, связанные с категориями. Предполагается, что каждая категория имеет относительно узкую область, например *нежирный творог*, поэтому продукты в пределах одной категории рассматриваются потребителем как взаимозаменяемые. Следовательно, несколько брендов могут конкурировать в одной категории за клиентов.

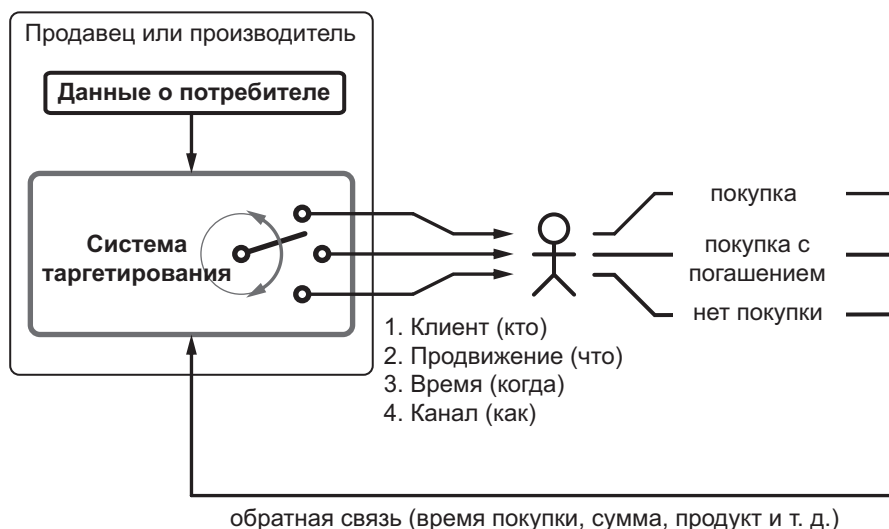


Рис. 3.1. Среда продвижения продаж в розничной торговле

- *Ретейлер* (розничный продавец) закупает товары у производителей и продает их потребителям с дополнительной наценкой. Ретейлер может также производить свою фирменную продукцию, именуемую *продуктом под частной маркой*, чтобы конкурировать с другими производителями в соответствующих категориях.
- *Целевое продвижение*, или *целевая рекламная акция* — это стимул, который можно применить к выбранным потребителям через маркетинговые каналы. Продвижение может заключаться в стимулировании с некоторым условием (например, *купи один — получи второй бесплатно*) или без условия, иметь денежное выражение в виде скидки или просто быть рекламой продукта или бренда. Продвижение может быть или не быть *погашаемым*, в том смысле, что потребителю может потребоваться представить доказательства участия в акции (штрих-код на печатном купоне или промокод), чтобы погасить часть стоимости покупки. Мы также используем слово *воздействие* как общий термин, который относится к продвижению и другим маркетинговым коммуникациям.
- Ретейлер владеет *маркетинговыми каналами*, такими как физические магазины или веб-сайты электронной коммерции, которые могут использоваться для информирования потребителей о рекламных акциях. Маркетинговые каналы нескольких ретейлеров могут объединяться в *сеть распространения рекламы*, которая управляется ретейлерами или сторонним агентством. Например, агентство может установить свои принтеры купонов в магазинах, принадлежащих нескольким розничным сетям.

Очень важно, чтобы ретейлер или агентство как владелец маркетингового канала могли получать информацию об отдельных потребителях и связывать воедино сделки, совершаемые одним и тем же потребителем или домохозяйством. Это часто реализуется с использованием *идентификаторов постоянных клиентов*, которые присваиваются клиентам посредством карт постоянного клиента или электронных учетных записей, идентификаторов кредитных карт или других сведений, доступных ретейлеру. Однако этот процесс часто несовершенен, и значительное число транзакций может оставаться анонимным.

- Рекламные акции могут распределяться по маркетинговым каналам в интересах как производителей, так и ретейлеров. Распределение может осуществляться как в пакетном режиме, когда электронные письма или печатные каталоги рассылаются большому числу клиентов, так и в режиме реального времени, когда рекламные акции генерируются в рамках отдельной транзакции, такой как покупка в магазине или посещение веб-сайта.
- Основные решения, которые должна принимать *система таргетирования* в отношении рекламных акций, заключаются в определении правильных по-

лучателей рекламы, правильных свойств рекламы, оптимального времени для ее предложения и правильного канала доставки.

- Предполагается, что ретейлер может идентифицировать потребителей, которым были направлены рекламные предложения, потребителей, купивших рекламируемый продукт, и, при необходимости, события погашения купонов рекламной акции. Обратите внимание, что покупка и погашение — это совершенно разные события, которые не следует путать: потребители, имеющие купон, не обязаны его погашать, и продукт, как правило, может быть куплен любым потребителем, хотя покупка может совершаться на разных условиях в соответствии с рекламными акциями. Помимо этих событий, система таргетирования может также получать дополнительные или внешние данные о потребителях, такие как демографические данные или ответы на опросы.

В среде, описанной выше, взаимодействие с потребителем часто структурировано в виде маркетинговых кампаний, которые являются удобной единицей оптимизации. Мы определяем *целевую кампанию* как маркетинговое действие, ограниченное бюджетом или продолжительностью, направленное на достижение определенной бизнес-цели путем распространения целевых предложений среди существующих клиентов или перспективных клиентов. Целевая кампания обычно включает следующие мероприятия.

ПЛАНИРОВАНИЕ. Планирование кампании обычно начинается с постановки бизнес-целей. На этапе планирования также должны быть определены основные свойства кампании, такие как бюджет, продолжительность или типы продвижения, которые могут быть выведены из цели.

ИСПОЛНЕНИЕ. Этап исполнения включает оценку потенциальных получателей и принятие решений о правильности предложений, сообщений, времени и каналов доставки.

ОЦЕНКА. Оценка показателей эффективности является критически важным мероприятием и может выполняться параллельно с исполнением для динамической корректировки.

Обратите внимание, что простой цикл, изображенный на рис. 3.1, не в полной мере отражает все важные аспекты управления продвижением. Во-первых, картина становится намного сложнее, когда речь заходит об управлении несколькими кампаниями или кампанией со сложной структурой, что часто случается на практике. Разные маркетинговые действия в рамках одной или нескольких кампаний могут влиять друг на друга, затрудняя принятие решений и оценку. Это означает, что порой недостаточно следить только за непосредственными событиями, свя-

занными с действием; иногда следует учитывать весь *жизненный цикл клиента*. Далее в этой главе мы обсудим, как программная система может справиться со всем этим. Второй важный фактор — отличие процесса таргетирования от подхода, изображенного на рис. 3.1, с точки зрения клиента. Жизненный цикл каждого отдельного клиента может включать несколько взаимодействий с ретейлером или производителем и, возможно, несколько каналов. Цепочка таких взаимодействий, называемая *циклом взаимодействия с клиентом*, должна обеспечивать непротиворечивый опыт во всех точках соприкосновения и на протяжении всего жизненного цикла. Она является ключевым аспектом в организации кампании, и мы подробно обсудим ее в разделе 3.6.1.

3.2. Бизнес-цели

Каждая маркетинговая кампания связана с определенными затратами и выгодами для каждого участника процесса, включая клиентов, ретейлеров, производителей и агентства. Концептуально каждая кампания должна иметь положительную *отдачу от инвестиций* — разницу между прибылью и затратами. Отдачу можно спрогнозировать до начала кампании или измерить после ее полного или частичного завершения. Предиктивные модели обычно оценивают отдачу по параметрам кампании, что позволяет оптимизировать экономику кампании.

Проблема, однако, в том, что доходы от кампании, как правило, имеют сложную структуру, включающую как денежные, так и неденежные компоненты, а также непосредственные и долгосрочные последствия. Эти последствия порой трудно измерить и еще труднее предсказать. В этом разделе мы обсудим некоторые основные соображения, касающиеся прибылей и убытков, а затем продолжим работу над более формальной структурой, которую можно использовать для моделирования кампаний. Эта структура оправдывает создание моделей таргетирования, как описано в следующем разделе, и закладывает основу для оптимизации кампании.

3.2.1. Производители и ретейлеры

Инициировать и финансировать маркетинговую кампанию может производитель или ретейлер. Во многих случаях они оба получают выгоду от увеличения продаж и числа лояльных клиентов. Однако особенности сотрудничества производителей и ретейлеров в значительной степени зависят от области бизнеса и маркетинговых стратегий для конкретных продуктов или их категорий. Детали этого сотрудничества играют важную роль для нас, потому что влияют на порядок предоставления или использования услуг программного таргетирования в розничной торговле.

Первое важное обстоятельство — стратегия управления взаимоотношениями с клиентами владельца клиентской базы, обычно ретейлера. Ретейлеры, торгующие товарами широкого потребления, как правило, приветствуют участие производителей в маркетинговом процессе, запрашивая *кампании, финансируемые производителем*. Такие кампании помогают производителям увеличить свою долю рынка в категории и выгодны ретейлеру. С другой стороны, элитные магазины, например модной одежды или косметики, позиционируют себя как персональные ассистенты и получают существенную добавочную стоимость от своих услуг по продаже. Такие ретейлеры не могут позволить третьим сторонам свободно взаимодействовать со своей клиентской базой. Поэтому они заранее закупают ассортимент товаров и управляют процессом их маркетинга, чтобы потом продать их с максимальной выгодой.

Второе обстоятельство — многие розничные торговцы предлагают продукцию под частной торговой маркой, что приводит к конфликту интересов с производителями. В случае рекламной услуги ретейлеры и производители могут договориться о специальных правилах, чтобы избежать деструктивной конкуренции в таких ситуациях, например путем исключения из таргетирования клиентов, очень лояльных к частной марке.

Наконец, ретейлеры заинтересованы в максимизации доходов в категории. Поощрение клиентов к переходу от высокодоходных продуктов к продуктам со скидками может быть вредным.

3.2.2. Затраты

Затраты на рекламную кампанию может нести как производитель, так и ретейлер. Но в любом случае и те и другие стремятся компенсировать расходы на кампанию более высоким объемом продаж. В мире товаров широкого потребления, например, кампании часто проводятся за счет производителя. Такие кампании инициируются производителем, а ретейлер фиксирует погашение купонов во время кампании и затем выставляет производителю счет с общей суммой затрат на погашение, которая обычно включает следующие компоненты.

РАСХОДЫ НА РАСПРОСТРАНЕНИЕ. Сюда входят затраты на дизайн купонов и их печать, отчисления маркетинговому агентству и фиксированные расходы, связанные с кампанией.

СТОИМОСТЬ ПОГАШЕНИЯ КУПОНОВ. Общая номинальная стоимость всех рекламных купонов. Ее можно оценить как произведение общего числа раздаваемых купонов, стоимости погашения одного купона и ожидаемой доли погашенных купонов.

ЗАТРАТЫ НА УСЛУГИ КЛИРИНГОВОЙ ОРГАНИЗАЦИИ. Купоны имеют свой жизненный цикл, предполагающий дополнительные расходы после их погашения. Когда покупатель вручает кассиру купон, тот кладет его в кассовый ящик или специальный конверт. В конце дня купоны складываются подобно наличным деньгам и упаковываются в мешки. Эти мешки затем отправляются в стороннюю клиринговую организацию. Служащие этой организации сортируют купоны, часто вручную, по накладной производителя и отправляют результаты подсчетов ретейлеру. Этот процесс влечет значительные расходы на услуги клиринговой организации, потому что крупный ретейлер может собирать миллионы купонов, а затраты на обработку одного купона порой сопоставимы с величиной скидки. Например, обработка одного купона клиринговой организацией в 2016 году стоила примерно 0,10 доллара США, тогда как скидка для большинства товаров широкого потребления по купону колебалась в диапазоне 0,50–2,00 доллара.

Структура затрат может меняться в зависимости от направления бизнеса и типа кампании, но обычно затраты легко поддаются оценке. С другой стороны, маркетинговые действия почти всегда связаны с некоторыми неденежными затратами или потерями, более сложными для оценки. В качестве простого примера можно привести так называемую *усталость от электронной почты* — уменьшение доли открываемых электронных писем и неудовлетворенность клиентов, вызванная слишком частыми или неактуальными электронными рассылками. Такие потери трудно оценить количественно, тем не менее далее мы увидим, что есть возможность сопоставить денежные показатели, такие как доходы, с соответствующими маркетинговыми действиями и тем самым дать количественную оценку и предсказать негативные последствия. Эти оценки потерь можно учесть в уравнении расчета прибыли.

3.2.3. Выгоды

Выгоды, связанные с кампанией, можно рассматривать с нескольких точек зрения. Самая заметная — увеличение объема продаж. Кампании, финансируемые производителями и ретейлерами, стимулируют потребителей делать покупки за счет затрат на кампанию, поэтому основное уравнение, описывающее выгоду от кампании, имеет следующий вид:

$$profit = Q(P - V) - C, \quad (3.1)$$

где Q — проданное количество, P — базовая цена за единицу, V — переменные маркетинговые затраты на единицу (средняя стоимость погашения купона, распространения и услуг клиринговой организации), а C — фиксированная стоимость рекламной кампании. Проще говоря, кампанию можно считать успешной, если

объем продаж Q_c , вызванный кампанией, превышает объем продаж Q_0 без кампании в степени, достаточной для покрытия затрат на кампанию:

$$Q_c(P - V) - C > Q_0 \cdot P. \quad (3.2)$$

Кампании, финансируемые производителем, направлены на достижение этой цели в контексте конкретного продукта, а также на увеличение доли производителя на рынке в соответствующей категории продуктов в долгосрочной перспективе. В то же время кампании, финансируемые производителем, как правило, также выгодны ретейлерам по следующим причинам.

- Рекламные акции стимулируют походы покупателя по магазинам. Производители и ретейлеры имеют общую цель — стимулировать больше походов по магазинам, поэтому рекламные кампании часто направлены на эту взаимовыгодную цель.
- Рекламные акции увеличивают размер корзины. Некоторые акции специально разработаны, чтобы заставить людей покупать больше данного продукта. Другие виды акций могут быть направлены на снижение расходов покупателя, чтобы высвободить деньги на дополнительные покупки.
- Рекламные акции повышают лояльность к ретейлеру. Потребитель воспринимает рекламные акции как результат сотрудничества между производителем и ретейлером, поэтому растет уровень доверия к ним обоим за их усилия по созданию добавленной стоимости и улучшению потребительского опыта.

Следовательно, ретейлер выигрывает от сотрудничества с производителем, которое влечет увеличение доходов и усиление эффекта лояльности. Именно по этой причине большинство продавцов товаров широкого потребления предоставляют производителям услуги по продвижению. Рекламные кампании, финансируемые ретейлерами, обычно направлены на продвижение частных марок, продвижение целых категорий продуктов или стимулирование оборота запасов. С точки зрения рекламы прибыль от кампаний, финансируемых ретейлерами, сравнима с прибылью от кампаний, финансируемых производителями, описанных выше. Однако перспективы оборачиваемости запасов разные, и мы обсудим это в главе 6, когда будем говорить об оптимизации цен и ассортимента.

Принцип максимизации объема продаж, приведенный в уравнении 3.2, является важным критерием для разработки рекламной кампании, но имейте в виду, что это очень упрощенный взгляд на управление отношениями с клиентами. Нам необходимо провести более тщательный анализ результатов кампании, чтобы лучше понять цели, которые можно использовать при разработке моделей таргетирования и кампаний. Поскольку кампания направлена на изменение отношений с потре-

бителем, ее цели лучше понять, изучая жизненный цикл клиента. Выделим три основных этапа взаимодействия между потребителем и брендом (производителем или ретейлером), следующих друг за другом, иногда многократно:

- Изначально потребитель не взаимодействует с брендом и предпочитает другие бренды или совершенно другие категории продуктов. Основная цель бренда на данном этапе — привлечь нового клиента.
- Клиентов, взаимодействующих с брендом, можно стимулировать покупать больше продуктов. Рекламные кампании для этих клиентов, как правило, следуют методологии увеличения продаваемого объема товара или комплексных продаж. Методология *увеличения продаваемого объема*, или *апселлинга* (up-selling), поощряет покупку большего объема товара, чем обычно покупает клиент у бренда. Методология *комплексных продаж* (cross-selling) стимулирует приобретение сопутствующих товаров.
- Наконец, клиент может перестать взаимодействовать с брендом. Обычно это называют *выбытием*, *бегством* или *оттоком клиентов*. Стоимость сохранения существующего клиента, как правило, намного меньше стоимости привлечения нового, поэтому бренд может делать специальные предложения клиентам, которые находятся на грани оттока.

Эти простые соображения образуют очень важную основу для управления взаимоотношениями с клиентами и продвижения товаров. Прежде всего отметим, что потребительское поведение и бизнес-цели сильно различаются на каждом из трех этапов жизненного цикла, как показано на рис. 3.2.



Рис. 3.2. Этапы жизненного цикла клиента

Потребители, находящиеся на первом этапе, должны привлекаться и переводиться в категорию постоянных клиентов с помощью маркетинговых действий,

специально разработанных для этой цели. На втором этапе клиентов необходимо стимулировать для максимального увеличения потребления. Наконец, необходимо своевременно выявлять и удерживать клиентов, находящихся на грани оттока. Эти три цели — *привлечение*, *максимизация* и *удержание* — являются очень популярной системой координат в маркетинге, которую можно использовать для ориентации отдельных кампаний и структурирования портфелей кампаний. Бренд должен уметь различать потребителей, находящихся на разных этапах, и это закладывает основу для процесса таргетирования. Как мы увидим далее, каждую из этих целей относительно легко можно отобразить в прогностическую модель, то есть данный набор целей таргетирования хорошо подходит для программирования.

С программной точки зрения разработка кампании на основе жизненного цикла требует решения двух основных задач. Первая — выявление потребителей, имеющих высокую *склонность* двигаться по кривой жизненного цикла, изображенной на рис. 3.3. При наличии возможности количественно оценить эту склонность мы сможем взаимодействовать с правильными потребителями, чтобы достичь цели и максимизировать выгоды.

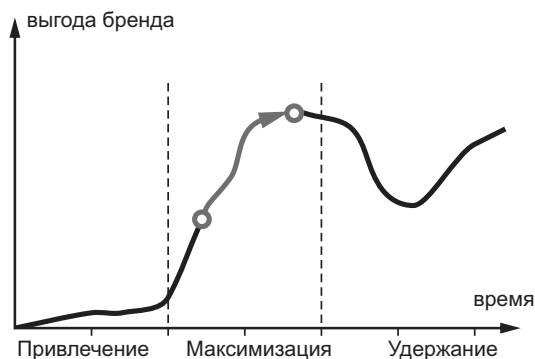


Рис. 3.3. Движение клиента по кривой жизненного цикла

Ориентация на правильных потребителей может повысить эффективность маркетинговых действий, но ее недостаточно для количественной оценки ожидаемой выгоды. Оценка ожидаемой выгоды — вторая по важности задача, требующая не только предсказать склонность потребителя перейти в определенную точку кривой жизненного цикла, но и оценить доход, который будет получен от потребителя после этой точки. Этот доход соответствует площади под кривой жизненного цикла. Однако измерять нужно не общую выгоду, а прибавку, которую дадут маркетинговые действия в сравнении со случаем, когда такие действия отсутствуют. Эта

прибавка соответствует заштрихованной области на рис. 3.4. Другими словами, модель кампании должна предсказывать прирост прибыли сверх доходов, которые дает стратегия бездействия, и ожидаемую выгоду от кампании, которая в свою очередь определяется с точки зрения склонности реагировать и потенциального воздействия. В следующем разделе мы разрабатываем более формальную структуру для этого типа моделирования и оценки.

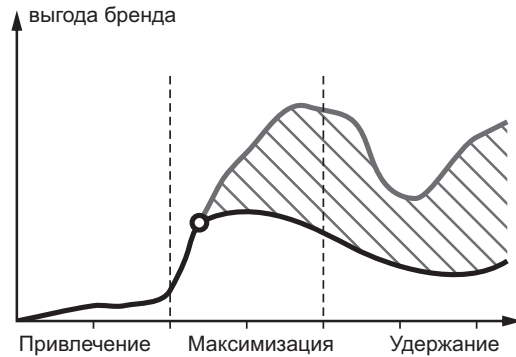


Рис. 3.4. Положительное влияние маркетингового действия. Верхняя кривая жизненного цикла соответствует периоду после маркетингового действия, а нижняя соответствует стратегии бездействия

3.3. Конвейер таргетирования

Определив среду и бизнес-цели, обсудим теперь подходы к задаче таргетирования и управления кампаниями в программной системе. Эту задачу можно рассматривать как создание процесса, получающего маркетинговый бюджет и бизнес-цели в параметрах, разбивающего их на кампании и выполняющего соответствующие маркетинговые действия. Этот процесс можно спроектировать по-разному, в зависимости от особенностей использования системы таргетирования, но концептуально схему процесса часто можно представить в виде конвейера, аналогичного изображенному на рис. 3.5.

Конвейер начинается с определения бюджета, который можно выделить для разных маркетинговых действий. На первом шаге определяется, как бюджет должен распределяться между возможными видами действий: каковы основные цели и как эти цели сбалансированы? Результатом этого шага является набор целей, таких как привлечение новых потребителей продукта А и удержание потребителей продукта В, а также параметры финансирования последующих этапов процесса. Второй шаг — эволюция каждой цели в маркетинговую кампанию, то есть

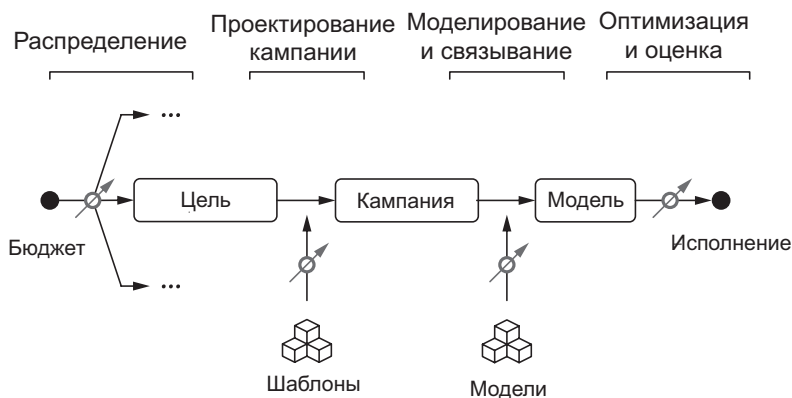


Рис. 3.5. Концептуальное представление конвейера таргетирования

проектирование кампаний. Программная система может использовать репозиторий *шаблонов кампаний*, которые можно выбирать и параметризовать на основе целей. Каждая кампания требует принятия ряда решений для определения целевых клиентов, оптимального времени таргетирования, параметров сообщения и т. д. Как правило, для этого используются предиктивные модели, которые необходимо обучить и связать с кампаниями. Модели возвращают оценки релевантности и другие сигналы, которые можно использовать для оптимизации параметров кампании, таких как список целевых клиентов или максимальная сумма скидки. Наконец, кампания исполняется, и собранные данные используются для дальнейшей оптимизации и оценки результатов. Программная система должна иметь возможность запустить конвейер как в режиме моделирования, для оценки разных стратегий, так и в режиме исполнения, для фактического таргетирования. Итак, конвейер таргетирования включает четыре основных элемента управления: распределение бюджета, проектирование кампаний, моделирование и оптимизацию исполнения. Мы обсудим эти элементы управления в следующих разделах, начав с моделей, которые служат основными строительными блоками, затем перейдем к проектированию и оптимизации кампаний и, наконец, рассмотрим общий бюджет и его распределение.

Конвейер таргетирования можно рассматривать не только с инженерной точки зрения, обозначенной на рис. 3.5, но и с точки зрения конечного пользователя (маркетолога). Эта перспектива очень важна, потому что описывает функции и свойства программной системы. Интерфейс системы сильно зависит от конкретных приложений и бизнес-среды, но для иллюстрации основных принципов можно рассмотреть простой гипотетический пример. Этот гипотетический поток управления кампанией показан на рис. 3.6.

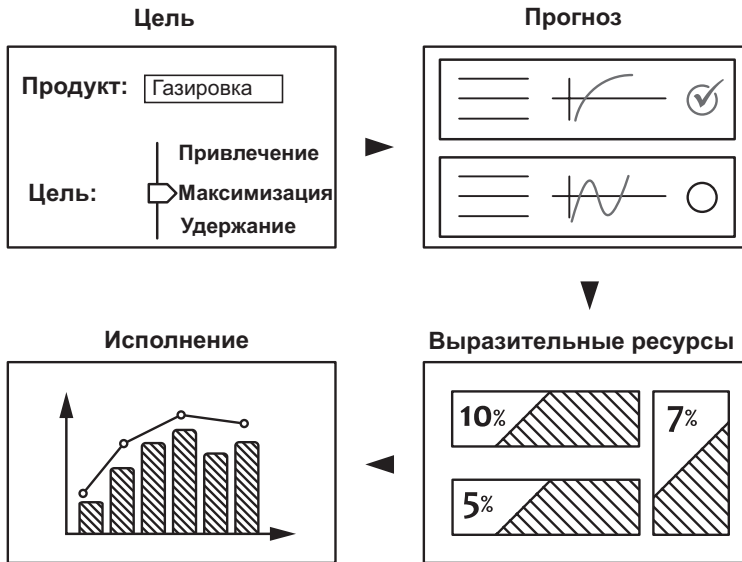


Рис. 3.6. Концептуальное представление процесса управления кампанией

Этот поток включает четыре основных этапа. При условии, что распределение бюджета выполняется заранее, точкой входа в процесс создания кампании является определение цели — выбор продвигаемых продуктов и конечных целей маркетинга. Для определения возможностей и предложения оптимальных стратегий система использует исторические данные, передовые приемы и предиктивные модели. Она прогнозирует ожидаемые результаты кампании, такие как затраты и прибыль, благодаря чему маркетолог может выбрать оптимальный вариант. После выбора шаблона кампании и оценки ее параметров система создает или настраивает выразительные ресурсы для взаимодействия с клиентами, такие как маркетинговые сообщения, шрифты и изображения. Совершенно понятно, что этот шаг требует участия человека. Наконец, сформированную кампанию можно запустить. Напомним еще раз, что это гипотетический и упрощенный поток, но он демонстрирует, чего мы пытаемся достичь в идеале при разработке системы таргетирования.

3.4. Моделирование и оценка отклика

Прежде чем углубиться в моделирование и разработку кампании, рассмотрим некоторые основные принципы моделирования и оценки отклика, чтобы получить представление о роли моделирования и оптимизации в системе таргетирования. Попутно введем новые понятия, такие как вероятность отклика, но пока не будем

указывать, как именно эти значения должны моделироваться и прогнозироваться. Оставим эти детали для последующего рассмотрения. Наша цель здесь — продемонстрировать, как затраты на кампанию, доходы и статистические свойства клиентов сочетаются в одной модели.

Цель рекламы и рекламных кампаний — изменить поведение потребителей и повлиять на их решения, побудить делать больше покупок, стимулировать покупку продвигаемых продуктов и т. д. Следовательно, успех кампании можно определить с точки зрения *отклика*, который можно оценить с помощью некоторых простых показателей, таких как процент погашенных рекламных купонов, или более сложных, включая прямые и косвенные, материальные и нематериальные выгоды. Эти показатели можно спрогнозировать до проведения кампании, чтобы помочь в принятии решений и их оптимизации или оценить постфактум на основе данных, собранных в ходе проведения. Эти две задачи одинаково важны, и мы обсудим их отдельно в последующих разделах, используя принципы моделирования жизненного цикла.

3.4.1. Платформа моделирования отклика

Платформа для моделирования отклика проста и универсальна. Она помогает разложить задачу моделирования кампании на несколько подзадач. Ее можно изменить и расширить с учетом сложности реальных маркетинговых кампаний. Начнем с относительно абстрактных условий, согласно которым бренду необходимо оптимизировать распределение рекламной акции или любого другого способа воздействия на потребителей, выбирая наиболее перспективных кандидатов для воздействия, чтобы максимизировать общий эффект кампании. Пока не будем уточнять, что понимается под эффектом — отложим рассмотрение этого аспекта до следующих разделов, — но предположим, что это какой-то количественный показатель, который можно сравнить с затратами. Задачи привлечения, максимизации и удержания можно считать вариантами этой задачи.

Напомним, что основная задача оптимизации маркетинга определяется как поиск стратегии, максимизирующей функцию эффекта. В случае отклика на кампанию моделируются общий эффект кампании с точки зрения вероятности отклика и ожидаемый чистый доход от клиента. Нашей целью оптимизации станет выбор клиентов, участвующих в продвижении, то есть аудитория кампании:

$$U_{opt} = \arg \max_{U \subseteq P} G(U), \quad (3.3)$$

где P — вся популяция потребителей, U — подмножество потребителей, достижимых в рамках кампании, и $G(U)$ ожидаемая прибыль от кампании, которая является

функцией стратегии таргетирования, выбирающей U из P . Ожидаемый доход от кампании можно смоделировать следующим образом:

$$G(U) = \sum_{u \in U} \Pr(R|u, T) \cdot (G(u|R) - C) + (1 - \Pr(R|u, T)) \cdot (-C), \quad (3.4)$$

где $\Pr(R|u, T)$ — вероятность отклика клиента u на воздействие (продвижение) T , $G(u|R)$ — чистая выгода от отклика клиента u , и C — стоимость ресурса продвижения. Первый член соответствует ожидаемой выгоде от отклика потребителя, а второй — ожидаемому убытку в случае отсутствия отклика на отправленное рекламное предложение. Цель состоит в том, чтобы максимизировать ожидаемую прибыль, подбирая подмножество клиентов, которые, вероятно, ответят наиболее выгодным способом. Уравнение 3.4 можно сократить, как показано ниже:

$$\begin{aligned} G(U) &= \sum_{u \in U} \Pr(\mathbf{R}|u, T) \cdot G(u|\mathbf{R}) - C = \\ &= \sum_{u \in U} \mathbb{E}[G|u, T] - C, \end{aligned} \quad (3.5)$$

где $\mathbb{E}[G|u, T]$ обозначает ожидаемую выгоду от отклика данного потребителя при условии, что тот получит рекламное предложение. Следовательно, критерий отбора клиентов можно упростить до

$$\mathbb{E}[G|u, T] > C, \quad (3.6)$$

потому что ожидаемый чистый доход неотрицателен и все потребители считаются независимыми. Далее оптимальное подмножество клиентов U можно определить как подмножество, максимизирующее эффект:

$$\operatorname{argmax}_{U \subseteq P} G(U) = \operatorname{argmax}_{U \subseteq P} \sum_{u \in U} \mathbb{E}[G|u, T] - C. \quad (3.7)$$

Обратите внимание, что этот подход можно интерпретировать как максимизацию чистого целевого дохода относительно случайного распределения ресурсов. Чтобы убедиться в этом, сравним эти два варианта с предположением, что в кампании примет участие фиксированное число клиентов $|U|$. Прирост дохода целевой кампании в сравнении с кампанией, распределяющей стимулы между $|U|$ клиентами, выбранными случайным образом, определяется как

$$\begin{aligned} &\operatorname{argmax}_{U \subseteq P} \sum_{u \in U} (\mathbb{E}[G|u, T] - C) - |U|(\mathbb{E}[G|T] - C) = \\ &= \operatorname{argmax}_{U \subseteq P} \sum_{u \in U} (\mathbb{E}[G|u, T] - \mathbb{E}[G|T]) = \\ &= \operatorname{argmax}_{U \subseteq P} \sum_{u \in U} \mathbb{E}[G|u, T], \end{aligned} \quad (3.8)$$

где $\mathbb{E}[G|T]$ — средний чистый доход на одного клиента в популяции. Этот средний чистый доход является константой, следовательно, его можно опустить, когда предполагается фиксированная мощность $|U|$. С другой стороны, тот же результат получится, если сократить уравнение 3.7, предположив, что $|U|$ является фиксированной величиной, а значит, можно снизить стоимость:

$$\operatorname{argmax}_{U \subseteq P} \sum_{u \in U} \mathbb{E}[G|u, T] - C = \operatorname{argmax}_{U \subseteq P} \sum_{u \in U} \mathbb{E}[G|u, T]. \quad (3.9)$$

Другими словами, случайный выбор получателей рекламы представляет собой базовый уровень, а задача максимизации стоимости эквивалентна перераспределению рекламных акций между группами потребителей.

Можно утверждать, что модель, определяемая уравнением 3.7, несовершенна, потому что отдаёт предпочтение потребителям, которые наверняка откликнутся на продвижение, но не учитывает клиентов, которые откликнутся в любом случае, принося ту же прибыль даже без продвижения [Radcliffe and Surry, 1999; Lo, 2002]. Следовательно, фактический доход от кампании по продвижению, по сравнению с базовым уровнем бездействия, может быть очень небольшим или даже отрицательным. Другой способ взглянуть на эту задачу — поставить следующий эксперимент: разделить группу клиентов, определяемую уравнением таргетирования 3.7, на две и воздействовать только на одну из них, а затем сравнить результаты. Может получиться так, что клиенты из первой группы будут активно погашать купоны и приобретать продукт, при этом клиенты из второй группы будут приобретать продукт в том же количестве или даже большем. Такая кампания явно неэффективна или даже вредна. Чтобы лучше понять проблему, рассмотрим отдельно следующие четыре возможные стратегии.

1. Выбрать группу клиентов $|U|$ согласно уравнению 3.7 и послать рекламные предложения каждому в этой группе.
2. Выбрать группу клиентов $|U|$ случайным образом и послать рекламные предложения каждому в этой группе.
3. Выбрать группу клиентов $|U|$ согласно уравнению 3.7, но не посылать рекламных предложений никому в этой группе.
4. Выбрать группу клиентов $|U|$ случайным образом, но не посылать рекламных предложений никому в этой группе.

Каждая из этих стратегий даст определенную прибыль для выбранной группы клиентов $|U|$, поэтому обозначим прибыль i -й стратегии как G_i . Уравнение 3.7 максимизирует разность $G_1 - G_2$, то есть эффективность таргетирования по сравнению со случайным распределением. Альтернативный подход, известный как *анализ*

дифференциального отклика, или моделирование увеличения эффективности, заключается в максимизации показателя эффективности, как показано ниже:

$$uplift = (G_1 - G_2) - (G_3 - G_4), \quad (3.10)$$

измеряющего не только увеличение эффективности в сравнении со случайным распределением, но также увеличение эффективности в сравнении с базовой стратегией бездействия для той же группы клиентов [Berry, 2009]. В этом случае уравнение 3.7 превращается в

$$\operatorname{argmax}_{U \subseteq P} \sum_{u \in U} \mathbb{E}[G | u, T] - \mathbb{E}[G | u, N] - c, \quad (3.11)$$

где второй член соответствует ожидаемой чистой выгоде для клиентов, которым не посылались рекламные предложения. Разницу между уравнениями 3.7 и 3.11 можно проиллюстрировать на примере следующей задачи: должен ли ретейлер предлагать скидку на чипсы тому, кто покупает их каждый день? В соответствии с уравнением 3.7 ответ на этот вопрос, скорее всего, будет положительным, потому что такой человек наверняка воспользуется купоном. При этом, вероятнее всего, клиент просто купит то же количество чипсов, но по более низкой цене, что фактически уменьшит прибыль ретейлера. Уравнение 3.11 устраняет эту проблему, учитывая поведение клиента по умолчанию. Обобщая пример, можно классифицировать клиентов по вероятности отклика на воздействие и вероятности отклика без воздействия, как показано на рис. 3.7.



Рис. 3.7. Классификация клиентов по эффективности воздействия

Анализ различий в вероятности позволяет предположить деление клиентов на четыре типа [Radcliffe and Simpson, 2007]: потерянные клиенты — с низкой вероятностью отклика на воздействия, — как правило, являются не лучшей целью для взаимодействия; надежные клиенты, откликающиеся независимо от воздействия, также выглядят не лучшей целью; клиентов, которых воздействие обычно отталкивает (их часто называют «не беспокоить»), тоже следует исключить из таргетирования; наконец, убеждаемые клиенты, которые наверняка откликнутся только при воздействии на них, — самые ценные. В следующих разделах мы подробнее рассмотрим моделирование увеличения эффективности и оптимизации прибыльности.

3.4.2. Оценка отклика

Платформа моделирования отклика — это базовый инструмент прогнозирования. Ее дополняет платформа оценки, которую можно использовать для оценки результатов кампании: действительно ли она помогла привлечь новых клиентов, заставила ли она существующих клиентов тратить больше или способствовала удержанию? Нам нужна возможность оценить эффективность с точки зрения отдачи инвестиций, которая определяется как дополнительный выигрыш от кампании по отношению к стратегии бездействия. Этот подход согласуется с принципами моделирования таргетирования и увеличения эффективности на основе жизненного цикла, изложенными в предыдущем разделе.

Стандартный подход к оценке дополнительных выгод заключается в сравнении двух групп потребителей: получивших рекламное предложение (*опытная группа*) и не получивших (*контрольная группа*). При наличии модели таргетирования обе группы обычно отбираются из числа клиентов с высокой вероятностью отклика, чтобы гарантировать статистическую согласованность.

Измеренное увеличение эффективности отражает влияние продвижения, независимо от стратегии таргетирования. Этот подход обычно реализуется путем исключения небольшого процента клиентов из целевой аудитории в самом конце процесса таргетирования, как показано на рис. 3.8.

Поведение групп сравнивается в течение некоторого периода времени, следующего за кампанией, обычно составляющего несколько циклов покупки данной категории продуктов, для получения стабильных результатов. Обратите внимание, что такой подход не требует следить за погашением купонов — нас интересует не доля погашенных купонов, а сравнение трат между двумя группами. Это очень удобная особенность, когда данные о погашении недоступны. Мы вернемся к статистическим деталям оценки в конце этой главы.

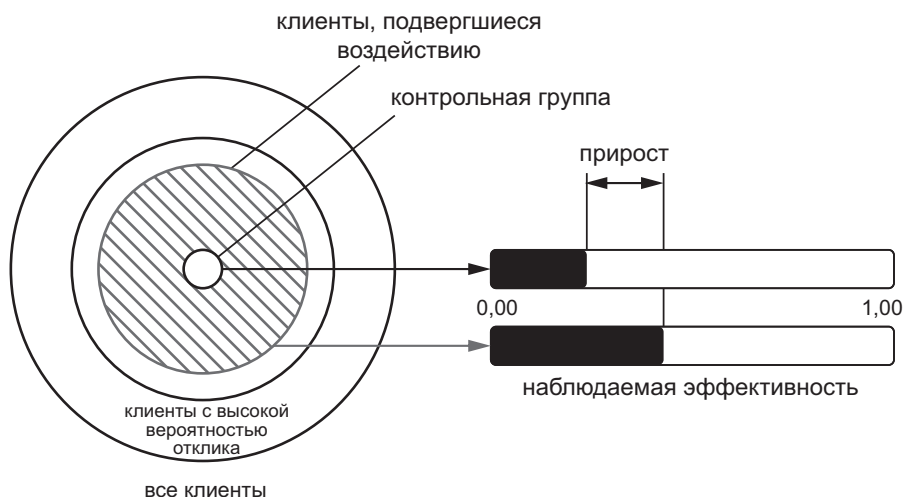


Рис. 3.8. Оценка эффективности продвижения делением клиентов на опытную и контрольную группы

3.5. Строительные блоки: таргетирование и модели ценности клиента

Модели таргетирования и пожизненной ценности клиента (lifetime value, LTV) — основные строительные блоки процесса таргетирования. Цель модели таргетирования — количественная оценка пригодности данного потребителя для конкретной бизнес-цели в данном контексте. Например, модель может оценить пригодность потребителя для рекламной кампании картофельных чипсов, учитывая, что рекламное предложение будет отправлено завтра по СМС. Модели могут создаваться для разных целей и контекстов, и система таргетирования часто поддерживает репозиторий моделей с сопутствующими метаданными, помогающими извлекать модели в соответствии с критериями. Например, в репозитории может иметься модель для кампании привлечения к категории картофельных чипсов и еще одна для кампании максимизации в категории газированных напитков. Модели — это базовые примитивы, которые можно объединять друг с другом и с другими строительными блоками для получения более сложных программных потоков. Из моделей можно собирать маркетинговые кампании, а из кампаний — маркетинговые портфолио.

Следует подчеркнуть, что программная система может использовать модели как для *прогнозирования*, так и для *управления*. Наиболее прямое применение — про-

гнозирование характеристик потребителей, таких как склонность к отклику на электронное письмо или ожидаемая пожизненная ценность. Многие модели, однако, относительно прозрачно выражают зависимость между входными и выходными параметрами, поэтому система может содержать дополнительную логику, использующую такое предписывающее понимание или, по крайней мере, способную дать маркетологу некоторые рекомендации. Например, параметры регрессионной модели, предсказывающей отклик, могут указывать на положительную или отрицательную корреляцию с определенными каналами коммуникации или другими параметрами, и это понимание можно использовать для дополнительных корректировок, таких как ограничение количества сообщений при наличии отрицательной корреляции. В этом разделе мы рассмотрим три основные категории моделей, которые можно использовать по отдельности или вместе:

МОДЕЛИ ПРЕДРАСПОЛОЖЕННОСТИ. Идея моделей предрасположенности состоит в оценке вероятности, что потребитель совершит определенное действие, например, купит определенный товар. Результатом таких моделей является оценка, пропорциональная вероятности, которую можно использовать для принятия решений о таргетировании.

МОДЕЛИ ВРЕМЕНИ ДО НАСТУПЛЕНИЯ СОБЫТИЯ. Модели предрасположенности оценивают вероятность события, но не оценивают наиболее вероятное время его наступления. Этот тип оценки играет важную роль во многих маркетинговых приложениях и требует использования другой статистической платформы.

МОДЕЛИ ПОЖИЗНЕННОЙ ЦЕННОСТИ. Модели LTV (Lifetime Value) используются для количественной оценки ценности клиента и влияния маркетинговых действий.

Начнем с обзора элементов и источников данных, используемых в моделировании, а затем обсудим некоторые традиционные методы, которые можно рассматривать как модели эвристической предрасположенности и оценки LTV. Эти методы обычно предполагают, что вероятность отклика и ценность клиента пропорциональны одной или нескольким основным характеристикам, таким как частота покупок. Затем они группируют клиентов в сегменты, чтобы весь сегмент можно было включить или исключить из определенной кампании. Эти методы можно рассматривать как таргетирование на основе правил. Затем мы определим более совершенные модели, применив статистические методы.

3.5.1. Сбор данных

Сбор и подготовка данных является одним из самых важных и сложных этапов моделирования. Подробное изучение методологии подготовки данных выходит за

рамки этой книги, тем не менее рассмотрим некоторые принципы, которые могут помочь упорядочить процесс и избежать ошибок при моделировании. Модели таргетирования и LTV обычно нацелены на прогнозирование потребительского поведения в зависимости от наблюдаемых показателей и свойств, поэтому важно собирать и использовать данные в соответствии с причинно-следственными зависимостями. С этой точки зрения элементы данных можно организовать по уровням, причем каждый уровень зависит от предыдущего:

ПЕРВИЧНЫЕ МОТИВАЦИИ. Поведение потребителя определяется такими базовыми факторами, как ценность продукта или услуги, вкусы, потребности, образ жизни и предпочтения. Многие из них нельзя наблюдать непосредственно, но некоторые сведения, такие как демографические данные или предпочитаемые маркетинговые каналы, можно получить с помощью регистрационных форм и опросов в рамках программы лояльности или приобрести у сторонних поставщиков данных.

ЭМПИРИЧЕСКИЕ МОТИВАЦИИ. Следующий уровень свойств создается взаимодействием между клиентом и брендом. Эти свойства характеризуют общее качество обслуживания клиентов, включая удовлетворенность, лояльность и структуру потребления. Некоторые из эмпирических свойств можно оценить количественно, прямо или косвенно с помощью таких показателей, как частота покупок.

ПОВЕДЕНИЕ. Наиболее важной категорией данных являются сведения о явно наблюдаемом поведении, такие как покупки, посещение веб-сайтов, история просмотров и переходы по ссылкам из электронных писем. Эти данные часто фиксируют взаимодействие с отдельными продуктами в определенные моменты времени. Поведенческие данные несут наиболее важные сигналы для моделирования.

РЕЗУЛЬТАТЫ. Наконец, действия клиентов напрямую влияют на финансовые показатели, такие как доход или прибыль. Важно помнить, что на самом деле эти показатели не объясняют движущие силы, определяющие поведение клиента; они просто фиксируют конечные результаты.

Данные, описанные выше, также должны быть связаны с дополнительными измерениями, такими как данные в каталоге, сезонность, цены, скидки и информация о магазине. Важно обеспечить возможность агрегирования данных на разных уровнях иерархических измерений для определения оптимального уровня детализации. Например, модель может использовать данные, агрегированные на уровне продукта, категории или отдела.

Процесс моделирования, как правило, должен ориентироваться не только на анализ результатов, но также на выявление скрытых свойств и причинно-следственных связей. Анализ финансовых результатов, конечно, важен, но часто желательно выявить также связи между маркетинговыми действиями и их результатами с при-

менением поведенческих концепций. Например, решение, выражающее доход как функцию от интенсивности рекламы, не всегда является достаточно информативным и действенным. Решение, количественно определяющее влияние рекламы на лояльность клиентов и поведенческие модели (например, переход из одного сегмента клиентов в другой), а затем связывающее свойства клиента с доходом, скорее всего, будет более информативным и действенным.

3.5.2. Многоуровневое моделирование

Модели таргетирования оценивают релевантность клиента для бизнес-цели на основе признаков, полученных из профиля клиента. Один из основных подходов заключается в использовании единственного показателя, например среднемесячной суммы, потраченной на бренд или категорию. Затем этот показатель можно использовать двумя способами. Во-первых, как показатель зависимости между клиентом и рекламной акцией, поскольку рекламные акции обычно проводятся для определенного бренда и категории. Следовательно, наиболее релевантные акции для данного клиента можно выбирать на основе бренда и категории с самыми высокими финансовыми показателями. Во-вторых, потребителей можно отсортировать по показателю и выбрать наиболее ценных из них для данной акции. Классическим примером такого подхода является многоуровневая сегментация, когда потребителям присваиваются золотой, серебряный и бронзовый уровни, в зависимости от их оценки и эвристически выбранного порога, как показано на рис. 3.9.

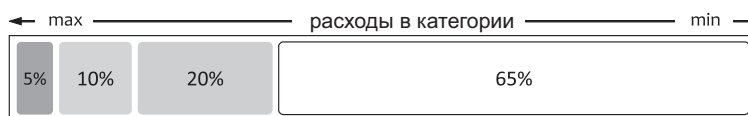


Рис. 3.9. Пример сегментирования по финансовому показателю (золотой — серебряный — бронзовый уровни). Клиенты сортируются по их расходам в категории в течение некоторого фиксированного периода времени. Верхним 5 % присваивается золотой уровень. Следующим 10 % — серебряный, и следующим 20 % — бронзовый. Остальные клиенты считаются неподходящими для рекламной акции

Каждому уровню приписываются такие показатели, как средняя ожидаемая частота откликов и средние расходы на одного клиента, рассчитанные на основе исторических данных. Для каждой рекламной акции можно определить оптимальное подмножество уровней, передавая затраты на продвижение и показатели уровня в платформу моделирования откликов. Например, можно определить, что одна кампания принесет прибыль, если будет нацелена только на золотой уровень, а другая даст максимальную прибыль при нацеливании на два уровня — золотой и серебряный.

Сегментацию по единственному показателю можно расширить, добавив дополнительные показатели. Как уже говорилось выше, жизненный цикл потребителя является важным фактором при разработке рекламных кампаний, поэтому способность ориентироваться на отдельные фазы жизненного цикла имеет большое значение. Этапы жизненного цикла характеризуются как общими расходами в категории, так и лояльностью к бренду, то есть относительными расходами на бренд по сравнению с другими брендами, поэтому клиентов можно классифицировать по сегментам с помощью этих двух показателей, как показано на рис. 3.10. Этот подход, известный как *сегментация по лояльности-расходам*, широко используется в традиционных кампаниях, финансируемых производителем.

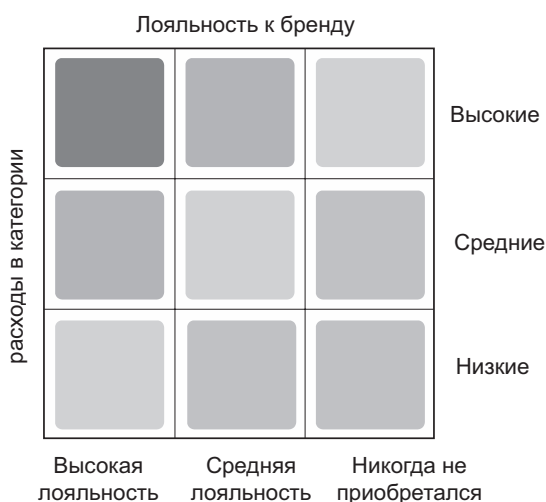


Рис. 3.10. Пример сегментации по лояльности-расходам

Клиенты, лояльные к бренду и тратящие много денег в категории, безусловно, являются самыми ценными и должны вознаграждаться и удерживаться. Клиенты с большими затратами в категории, но не лояльные к данному бренду, являются лучшими кандидатами на пробные предложения, и так далее. По аналогии с многоуровневой сегментацией рекламные акции можно нацеливать на оптимальное подмножество сегментов, начиная с верхнего левого угла сетки на рис. 3.10 и оценивая потенциальный результат включения дополнительных сегментов, пока не будет достигнут нижний правый угол. Такой подход можно рассматривать как упрощенный метод таргетирования, прогнозирующий потребительскую ценность на основе двух показателей — расходов в категории и доли бренда в кошельке. Это очень грубые критерии, которые можно уточнить с помощью методов предиктивного моделирования.

3.5.3. RFM-моделирование

Другой популярный подход к сегментации — так называемый анализ по *давности-частоте-расходам* (Recency-Frequency-Monetary, RFM). Он напоминает подход на основе лояльности-расходов, но использует три показателя:

ДАВНОСТЬ. Количество единиц времени, прошедших с момента последней покупки, сделанной клиентом. Этот показатель может измеряться непосредственно в единицах времени (например, месяцах) или отображаться в некоторой оценке. Например, клиентов можно отсортировать по возрастанию давности последней покупки, отобрать первые 20 % и присвоить им оценку 5, затем отобрать следующие 20 % и присвоить им оценку 4, и так далее, до последних 20 % с оценкой 1.

ЧАСТОТА. Среднее количество покупок в единицу времени. И снова показатель может измеряться в непосредственных единицах или оценках.

РАСХОДЫ. Общая сумма, потраченная в единицу времени. Этот показатель обычно измеряется с использованием интервалов или оценок.

В большинстве случаев для всех трех показателей используется одна и та же дискретная шкала оценок, скажем, от 1 до 5. В этом случае модель RFM можно рассматривать как трехмерный куб, состоящий из ячеек, каждая из которых определяется тройкой значений показателей и соответствует клиентским сегментам. Решения о таргетировании можно принимать, выбирая подмножество сегментов из куба RFM. Один из возможных методов — суммирование всех трех показателей в один балл и выбор клиентов, чей балл превышает пороговое значение, — это соответствует срезанию угла куба RFM.

Анализ RFM основан на эмпирическом наблюдении высокой корреляции показателей давности, частоты и расходов с вероятностью отклика и пожизненной ценностью. Несмотря на обоснованность этого предположения, RFM-подход является поверхностным, поскольку оценивает конечный результат маркетинговых процессов и действий потребителей, а не факторы, влияющие на потребительское поведение. Как будет показано в следующем разделе, с помощью кластеризации можно получить более гибкое решение.

3.5.4. Моделирование предрасположенности

Простые модели сегментирования и RFM-анализ можно рассматривать как частный случай регрессионного анализа с очень ограниченным набором признаков и эвристических допущений относительно взаимосвязей между показателями и ожидаемыми результатами. Следующим нашим шагом будет создание более формальной модели оценки.

Цель моделирования predisposedness — поиск клиентов с относительно высокой вероятностью определенного поведения или выполнения определенных действий в будущем. Количество действий, которые можно спрогнозировать и использовать в таргетировании, очень велико. Рассмотрим несколько типичных примеров:

ПРЕДРАСПОЛОЖЕННОСТЬ ПРОБОВАТЬ НОВЫЕ ПРОДУКТЫ. Потребители, в настоящее время не покупающие определенный продукт, но имеющие высокую predisposedness купить его в будущем, являются хорошей целью для кампаний по привлечению.

ПРЕДРАСПОЛОЖЕННОСТЬ К РАСШИРЕНИЮ КАТЕГОРИИ. Потребители, имеющие высокую predisposedness переходить из одной категории продуктов в другую или пробовать новую категорию, являются хорошей целью для кампаний по увеличению продаж или продаж сопутствующих товаров. Примерами такой аудитории могут служить потребители, которые, скорее всего, откажутся от покупки обычных продуктов в пользу продуктов премиум-класса.

ПРЕДРАСПОЛОЖЕННОСТЬ К БОЛЬШОМУ ЧИСЛУ ПОКУПОК. Потребители, которые наверняка готовы увеличить среднее количество покупок продукта, являются хорошей целью для кампаний, направленных на максимизацию доходов.

ПРЕДРАСПОЛОЖЕННОСТЬ К ОТТОКУ. Клиенты, которые могут отказаться от подписки на услугу или прекратить покупку продукта, являются хорошей целью для кампаний, нацеленных на удержание.

ПРЕДРАСПОЛОЖЕННОСТЬ К ВОВЛЕЧЕНИЮ. Предрасположенность к вовлечению — это вероятность отклика на маркетинговое действие, например на ссылку в электронном письме.

ПРЕДРАСПОЛОЖЕННОСТЬ К ИЗМЕНЕНИЮ ПОКУПАТЕЛЬСКИХ ПРИВЫЧЕК. Каждый клиент имеет покупательские привычки, которые в конечном счете определяют его ценность: как часто он совершает покупки, какие продукты покупает и в из каких категорий, и т. д. Обычно эти привычки не меняются с течением времени, и после того как бренд сумеет изменить уровень взаимодействия с клиентом, этот уровень имеет тенденцию оставаться постоянным. То есть в общем случае бренды заинтересованы в поиске клиентов, которые predisposedness менять свои привычки. Например, людей, переехавших в другой город, окончивших школу или университет, только что женившихся или вышедших замуж и т. д. Классическим примером такого моделирования является попытка компании Target¹

¹ Американская компания, управляющая сетью магазинов розничной торговли, работающих под марками Target и SuperTarget. — *Примеч. пер.*

спрогнозировать беременность клиентов на ранних стадиях, потому что роды, очевидно, меняют покупательские привычки [Duhigg, 2012].

Обратите внимание, что основные маркетинговые цели — привлечение, максимизация и удержание — можно выразить на языке предрасположенности. Подход, основанный на предрасположенности, удобен с точки зрения моделирования откликов, потому что позволяет оценить рентабельность кампании путем умножения ожидаемых прибылей и убытков на прогнозируемые вероятности результатов.

3.5.4.1. Моделирование методом аналогии

Метод аналогий — один из наиболее важных методов моделирования предрасположенности. Он основан на наблюдении, что предрасположенность — это фактически вероятность перехода клиента из одной точки на кривой жизненного цикла в другую, благодаря чему можно обучить предиктивную модель, используя профили потребителей, проявивших такое поведение в прошлом, и затем использовать обученную модель для оценки предрасположенности данного клиента по его профилю. Например, профили клиентов, которые какое-то время не покупали данный продукт, а затем начали его покупать, можно использовать для обучения модели, идентифицирующей клиентов с высокой предрасположенностью впервые попробовать этот продукт.

Метод аналогий относится к категории задач классификации, поэтому для него необходимо отобрать профили и задать метки отклика. Предполагается, что профиль клиента может включать индивидуальные атрибуты, такие как доход или размер домохозяйства, а также коллекцию поведенческих событий, каждое из которых сопровождается отметкой времени. Все события для каждого профиля помещаются на временную линию, и определяются три последовательных периода: период наблюдения, промежуточный период и период итогов. Эти периоды показаны на рис. 3.11. Период наблюдения используется для создания признаков, а период итогов — для создания меток откликов. Эти два периода может разделять промежуточный период, который используется, если нужно предсказывать события в относительно отдаленном будущем. Например, модель, которая предсказывает отток клиентов, вероятно, следует обучать с периодом итогов, смещенным в будущее, — бессмысленно предсказывать немедленный отток клиентов, потому что это не дает времени для выполнения каких-либо смягчающих маркетинговых действий.

Модель обучается на наборе исторических профилей, содержащих как интервалы наблюдения, так и интервалы итогов. Затем эта модель используется для оценки текущих профилей, которые, конечно же, содержат только наблюдаемую часть, и прогнозирования ожидаемого результата.

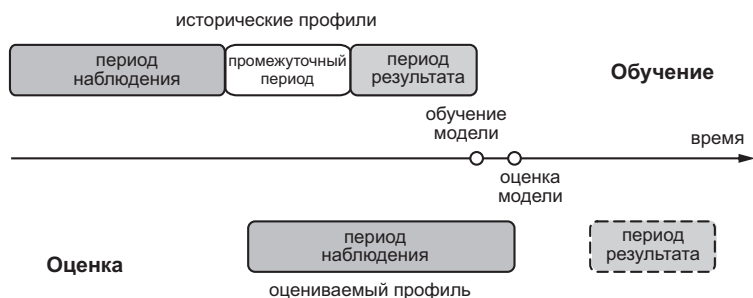


Рис. 3.11. Временные периоды моделирования методом аналогий

Проектирование признаков и меток отклика в значительной степени зависит от конкретной области. В этом разделе мы сосредоточимся на розничной торговле, а проектирование признаков для онлайн-рекламы рассмотрим в следующих разделах. В методе аналогий, как правило, можно использовать любые данные клиента, включая демографию, покупки и отклик на маркетинговые действия, такие как переходы по ссылкам в электронных письмах и погашение купонов. Признаки обычно указываются в виде различных комбинаций временных периодов, показателей и фильтров, которые могут применяться к данным в профилях. Такой процесс для данных о покупках показан на рис. 3.12. Во-первых, признаки могут вычисляться для разных периодов в пределах периода наблюдения. Эти подпериоды обычно откладываются от конца интервала наблюдения: прошлый месяц, последние три месяца, последние шесть месяцев и т. д. В пределах подпериода можно рассчитать разные показатели, такие как сумма расходов или частота покупок, и применить различные фильтры, например по категории, бренду, продукту, типу оплаты или дням недели. Наконец, ценность можно выразить в таких единицах, как рубли или дни, проценты, средняя стоимость чека или в виде бинарных переменных да/нет. Например, кривая, соединяющая прямоугольники на рис. 3.12, соответствует доле категории хлебобулочных изделий в потребительских расходах за последние 6 месяцев, по отношению к другим категориям. Такой подход позволяет получить относительно большое количество признаков, которые можно использовать для обучения и оценки предиктивной модели. Такой же подход можно использовать для данных об откликах на маркетинговые действия и данных из цифровых каналов.

Метка отклика генерируется из периода итогов в соответствии с целью. Например, если модель создана для прогнозирования клиентов с высокой predisposedностью опробовать данный продукт, то в метке ответа будет указано, был ли приобретен этот продукт. Другой пример — кампания удержания, где метка отклика указывает, ушел ли клиент. Набор обучающих профилей можно также предварительно отфильтровать согласно задаче. В примере с predisposedностью

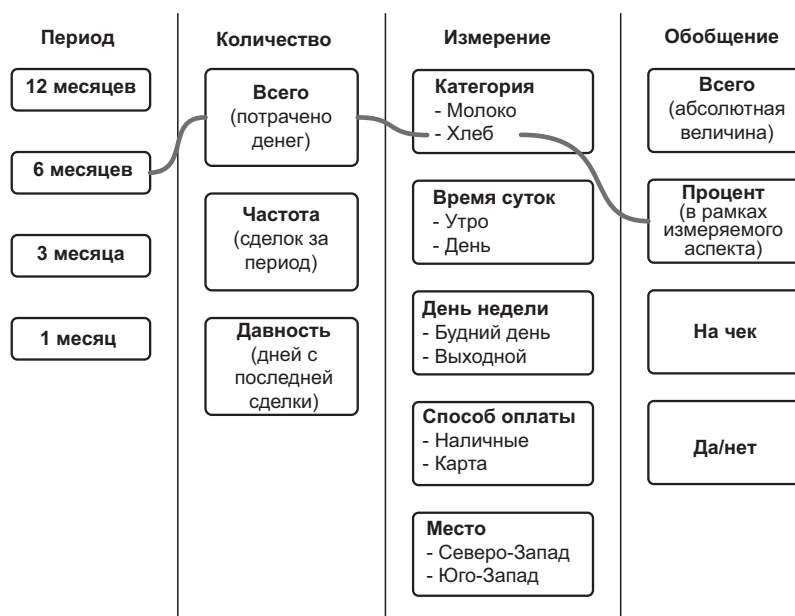


Рис. 3.12. Пример проектирования признаков профиля для данных о покупках

пробовать новые продукты в обучающий набор должны входить только клиенты, не покупавшие продукт в период наблюдения (а затем начавшие его покупать в период итогов — таких клиентов называют *естественными любителями пробовать новое*). То же правило применяется для оценки модели — клиент, который уже покупает продвигаемый продукт, сразу исключается как неподходящий.

ПРИМЕР 3.1

Вот небольшой пример, иллюстрирующий моделирование методом подобию. Рассмотрим следующий сценарий: ретейлер, торгующий молочными и хлебобулочными изделиями, шесть месяцев назад начал работать с новым брендом — производителем молочных десертов, — и теперь бренд попросил запустить кампанию по привлечению новых потребителей, финансируемую производителем. Целью таргетирования в данном случае является выявление клиентов с высокой предрасположенностью попробовать десерт. Предположим, что система таргетирования сформировала набор обучающих данных с 12 историческими профилями клиентов и 5 признаками, как показано в табл. 3.1.

Таблица 3.1. Набор обучающих данных для моделирования методом подоби

№ п/п	Всего хлебобулочных	Хлебобулочных в выходные	Всего молочных	Молочных в выходные	Оплата картой	Отклик
1	150	10	150	140	1	1
2	210	20	120	110	1	0
3	190	190	210	20	1	1
4	270	250	190	0	1	1
5	180	180	190	10	1	1
6	260	250	230	20	0	1
7	270	30	210	210	1	0
8	150	40	150	50	1	0
9	90	70	120	100	0	0
10	30	0	200	200	1	0
11	190	190	250	10	1	1
12	10	0	30	0	1	0

«Всего хлебобулочных» и «Всего молочных» — это общие расходы в соответствующих категориях за период наблюдения. «Хлебобулочных в выходные» и «Молочных в выходные» — это расходы в тех же категориях в выходные дни недели, то есть сумма, потраченная в рабочие дни, равна разнице сумм «Всего» и «В выходные» в каждой категории. Столбец «Оплата картой» определяет способ оплаты — картой или наличными. Наконец, переменная «Отклик» указывает, стал ли клиент покупать десерт. После беглого обзора этого небольшого набора данных можно сделать вывод, что естественные любители пробовать десерт — это в основном клиенты, покупающие много хлебобулочных изделий в выходные дни и много молочных продуктов в рабочие. Для построения модели мы решили использовать логистическую регрессию, хотя на практике часто используются и другие варианты, включая деревья решений, случайные леса и наивную байесовскую модель. Обучив модель логистической регрессии, получаем оценки параметров, представленные в табл. 3.2.

Таблица 3.2. Параметры логистической функции для обучающего набора в табл. 3.1

Параметр	Оценка
Всего хлебобулочных	0,0012
Хлебобулочных в выходные	0,0199
Всего молочных	−0,0043
Молочных в выходные	−0,0089
Оплата картой	−0,4015

Обратите внимание, что расходы на хлебобулочные изделия положительно коррелируют с предрасположенностью попробовать продукт, тогда как расходы на молочные продукты имеют отрицательную корреляцию¹. Применяв модель к шести профилям с различными пропорциями расходов на хлебобулочные и молочные продукты, получим оценки коэффициента склонности, показанные в табл. 3.3. Как видите, высокую предрасположенность попробовать продукт имеют клиенты с высокими расходами на хлебобулочные изделия и низкими — на молочные продукты, независимо от способа оплаты. В реальной жизни такой результат можно интерпретировать, например, так: клиенты, активно покупающие обе категории товаров, считают молочные десерты заменителями хлебобулочных десертов.

Таблица 3.3. Прогнозируемая склонность попробовать продукт

Всего хлебобулочных	Хлебобулочных в выходные	Всего молочных	Молочных в выходные	Оплата картой	Склонность попробовать
10	0	50	50	1	0,26
20	20	200	200	1	0,07
150	20	100	30	1	0,37
250	20	190	30	1	0,31
250	200	190	30	1	0,94
250	200	190	30	0	0,96

¹ Подробное объяснение логистической регрессии вы найдете в главе 2. В целях упрощения в этом примере мы опустили типичные этапы, такие как проверка и диагностика модели.

Обратите внимание, что в этом примере мы не используем исторические отклики в качестве признаков, то есть не учитываем, откликнулся ли клиент на рекламные акции в прошлом. В реальной жизни это важный сигнал, помогающий повысить точность таргетирования, хотя вполне можно создавать модели без признака, описывающего отклик, если эти данные недоступны.

3.5.4.2. Моделирование отклика и увеличения расходов

Самые простые модели подобия, аналогичные описанным в предыдущем разделе, оценивают безусловную вероятность отклика на определенное действие. Включив маркетинговые коммуникации в набор функций, можно создать модель предрасположенности, оценивающую условную вероятность отклика (действия) на заданное маркетинговое воздействие. Одна из методологий создания таких моделей называется *пилотные кампании*. Идея состоит в том, чтобы первоначально провести рекламную акцию для относительно небольшой группы получателей, собрать отклики и создать модель классификации, максимизирующую разницу между откликнувшимися и неоткликнувшими респондентами. Это соответствует модели подобия, обученной на совокупности профилей клиентов, подвергнувшихся воздействию, с использованием индикатора отклика в качестве метки обучения. Эта модель оценивает вероятность отклика на воздействие как

$$\Pr(R|T, \mathbf{x}), \quad (3.12)$$

где R — индикатор отклика, T — индикатор воздействия, а \mathbf{x} — вектор признаков профиля. Созданную модель можно использовать для проведения полномасштабной кампании, то есть для таргетирования клиентов с высокой предрасположенностью реагировать на воздействие. Иногда есть возможность создать модель для аналогичных кампаний с использованием исторических данных, без запуска пилотной кампании. Таким образом, традиционные модели предрасположенности предназначены для выявления клиентов, склонных откликнуться на рекламные акции или другие маркетинговые коммуникации. Недостатком такого подхода является то, что такие модели не отличают клиентов, которые откликнутся в любом случае, даже без воздействия на них. Другими словами, модель предрасположенности предсказывает высокую вероятность отклика, но после проведения кампании наблюдаемая разница между опытной и контрольной группами, то есть увеличение эффективности, может оказаться незначительной, более того, контрольная группа может даже превзойти тестовую. Мы уже затронули эту проблему в контексте платформы моделирования отклика, но теперь нам нужно копнуть глубже и определить, как решить ее в моделировании предрасположенности.

Проблема с увеличением эффективности возникает из-за того, что описанный выше процесс моделирования предрасположенности учитывает только клиентов, подвергшихся воздействию. Это делает структурно невозможным моделирование увеличения эффективности. Проблему можно обойти, добавив контрольную группу в пилотную кампанию. Эта группа должна включать случайно выбранные профили клиентов, не участвовавших в пилотной кампании. В этом случае можно наблюдать результаты для четырех групп: клиенты, подвергшиеся воздействию и откликнувшиеся, подвергшиеся воздействию и неоткликнувшиеся, и клиенты из контрольной группы, откликнувшиеся и неоткликнувшиеся, как показано на рис. 3.13.

Воздействие	Нет	Откликнувшиеся из контрольной группы (CR)	Неоткликнувшиеся из контрольной группы (CN)
	Да	Откликнувшиеся участники кампании (TR)	Неоткликнувшие участники кампании (TN)
		Да	Нет
		Отклик	

Рис. 3.13. Оценка групп в моделировании предрасположенности по эффективности [Kane et al., 2014]

Наличие четырех наблюдаемых групп позволяет создать модель, максимизирующую увеличение эффективности, то есть разницу между частотой отклика в опытной и контрольной группах:

$$uplift(\mathbf{x}) = \Pr(R|T, \mathbf{x}) - \Pr(R|C, \mathbf{x}), \quad (3.13)$$

где первый член — вероятность отклика после воздействия, а второй — вероятность отклика индивида из контрольной группы. Эти две вероятности можно оценить с помощью двух отдельных моделей классификации, обученных на опытной и контрольной группах соответственно или с помощью одной модели, обученной на объединении опытных и контрольных профилей с добавлением дополнительного признака воздействия [Lo, 2002]. Проблема с подходом на основе двух моделей заключается в том, что отдельно созданные модели могут иметь несопоставимые

шкалы оценки и выбирать признаки, которые фактически не являются прогнозирующими в отношении увеличения эффективности, поэтому это решение зачастую не позволяет добиться лучших результатов, чем базовая модель предрасположенности [Radcliffe and Surry, 2011; Kane et al., 2014]. Подход на основе одной модели позволяет добиться лучших результатов, но может потребовать создания более сложной модели. Например, если в качестве базового метода моделирования использовать логистическую регрессию, вектор признаков должен включать не только признаки из профилей, но также индикаторы воздействия, чтобы модель имела функциональную форму

$$f(\mathbf{x}, \mathbb{I}(T) \cdot \mathbf{x}, \mathbb{I}(T)), \quad (3.14)$$

где $\mathbb{I}(T)$ — функция-индикатор, равная единице, если клиент \mathbf{x} был подвергнут воздействию, и нулю в противном случае [Lo, 2002]. Следовательно, увеличение эффективности оценивается как

$$uplift(\mathbf{x}) = f(\mathbf{x}, \mathbf{x}, 1) - f(\mathbf{x}, \mathbf{0}, 0). \quad (3.15)$$

Можно утверждать, что еще более точные результаты дает полиномиальная модель, предсказывающая вероятности для каждого из квадрантов на рис. 3.13 [Kane et al., 2014]. Ниже показано, как создать такую модель, чтобы выразить увеличение эффективности:

$$\begin{aligned} uplift(\mathbf{x}) &= \Pr(R|T, \mathbf{x}) - \Pr(R|C, \mathbf{x}) = \\ &= \Pr(R|T, \mathbf{x}) - (1 - \Pr(N|C, \mathbf{x})) = \\ &= \Pr(R|T, \mathbf{x}) - \Pr(N|C, \mathbf{x}) - 1, \end{aligned} \quad (3.16)$$

где N обозначает случай отсутствия отклика. Используя правило Байеса и тот факт, что $\Pr(T|\mathbf{x}) = \Pr(T)$, потому что опытная и контрольная группы выбираются случайным образом, получаем:

$$\begin{aligned} uplift(\mathbf{x}) &= \frac{\Pr(TR|\mathbf{x})}{\Pr(T|\mathbf{x})} + \frac{\Pr(CN|\mathbf{x})}{\Pr(C|\mathbf{x})} - 1 = \\ &= \frac{\Pr(TR|\mathbf{x})}{\Pr(T)} + \frac{\Pr(CN|\mathbf{x})}{\Pr(C)}. \end{aligned} \quad (3.17)$$

Применив те же преобразования к вероятности отклика в опытной группе, увеличение эффективности можно выразить так же как:

$$\begin{aligned} uplift(\mathbf{x}) &= (1 - \Pr(N | T, \mathbf{x})) - \Pr(R | C, \mathbf{x}) = \\ &= 1 - \frac{\Pr(TN | \mathbf{x})}{\Pr(T)} - \frac{\Pr(CR | \mathbf{x})}{\Pr(C)}. \end{aligned} \quad (3.18)$$

Сложив уравнения 3.17 и 3.18, получаем окончательное выражение оценки увеличения эффективности:

$$\begin{aligned} 2 \cdot uplift(\mathbf{x}) &= \frac{\Pr(TR | \mathbf{x})}{\Pr(T)} + \frac{\Pr(CN | \mathbf{x})}{\Pr(C)} - \\ &\quad - \frac{\Pr(TN | \mathbf{x})}{\Pr(T)} - \frac{\Pr(CR | \mathbf{x})}{\Pr(C)}, \end{aligned} \quad (3.19)$$

где вероятности в числителях оцениваются с помощью единой регрессионной модели. Оценку увеличения эффективности часто можно использовать в качестве альтернативы вероятности отклика, оцениваемой по базовым моделям предрасположенности. Как будет показано позже, система таргетирования может оптимизировать рентабельность кампании, выбирая получателей рекламы с самым высоким показателем эффективности вместо показателя предрасположенности.

3.5.5. Сегментирование и персонифицированное моделирование

Поведенческое сегментирование — это процесс разделения клиентов на группы, или *сегменты*, таким образом, чтобы клиенты в одном сегменте были похожи друг на друга, но отличались от клиентов в других сегментах. С точки зрения маркетинга сегментация обычно считается одним из самых важных, ценных, информативных и сложных проектов.

Целью этого процесса, как правило, является определение небольшого числа хорошо дифференцированных сегментов с четким смысловым значением, которые можно использовать для принятия стратегических решений. Результаты процесса сегментации обычно включают профили и модели сегментов, также называемые моделями кластеризации. Профиль сегмента включает отличительные свойства и показатели, а также некоторую интерпретацию типичного *образа* клиента. Упрощенный пример профилей сегментов приводится в табл. 3.4. Набор отличительных свойств обычно определяется применением алгоритмов кластеризации к набору исторических профилей клиентов, поэтому каждый сегмент соответствует группе существующих клиентов, а профиль сегмента представляет набор статистических показателей для этой группы. Изначально сегмент является простым списком су-

ществующих клиентов, однако его можно преобразовать в модель кластеризации, по сути, правило классификации для сопоставления любого профиля клиента с образом. Представление сегмента на основе модели играет важную роль, потому что может динамически относить клиентов к сегментам, в зависимости от признаков в их профилях.

Таблица 3.4. Пример сегментов и их показателей. Каждый сегмент можно интерпретировать в психографическом и поведенческом смыслах. Например, лица, ценящие удобства, очевидно менее чувствительны к ценам и имеют меньше детей, чем потребители из других сегментов. Этот сегмент содержит относительно небольшое количество клиентов, но вносит большой вклад в доход

Типичный образ	Сегмент 1	Сегмент 2	Сегмент 3
	ценящие удобства	случайные покупатели	любители выгодных покупок
Процент рынка	20	50	30
Процент доходов	40	40	20
Доля одежды	40	60	60
Доля электроники	50	20	10
Доля игрушек	10	20	30
Доля погашенных купонов	0,02	0,05	0,08

Обратите внимание, что поведенческое сегментирование сильно отличается от RFM-анализа, даже при том что последний можно рассматривать как разновидность сегментирования. RFM-анализ сегментирует клиентов по наблюдаемым финансовым результатам, тогда как поведенческое сегментирование направлено на выявление признаков, вызывающих этот результат. Во многих случаях признаки результата, такие как расходы, намеренно исключаются из признаков профилей перед кластеризацией, чтобы гарантировать создание сегментов на основе поведенческих особенностей, а не финансовых результатов. Второе важное отличие состоит в том, что RFM-анализ и его вариации используют фиксированный набор признаков, тогда как сегментирование является методом выявления наиболее характерных признаков. Эти свойства поведенческого сегментирования очень важны для стратегической маркетинговой аналитики, поскольку помогают понять силы, управляющие поведением клиентов (например, почему в одном сегменте наблюдается больший отток клиентов, чем в другом), и дифференцировать маркетинговые стратегии для каждого сегмента, используя отличительные характеристики. Это различие можно распространить, например, на выделение менеджеров для каждого сегмента.

Программная перспектива сегментирования для маркетинговой аналитики отличается от той, которую мы только что рассмотрели, потому что программная организация больше ориентирована на исполнение и тактические аспекты, а не на стратегию. Программная система таргетирования чаще оказывается пользователем результатов, полученных в процессе поведенческого сегментирования. Во-первых, признаки типичного представителя часто используются в качестве признаков в моделировании методом подобия и других правилах и моделях таргетирования, и неважно, как именно создаются эти признаки. Признаки типичного представителя несут важный сигнал о потребительском поведении и, следовательно, могут иметь значительную предсказательную силу для моделирования предрасположенности. Второе важное применение результатов сегментирования — моделирование на уровне сегмента. Модели предрасположенности, созданные для всей совокупности клиентов, часто имеют ограниченную точность, потому что предрасположенности могут определяться разными факторами. Например, отток клиентов в одном сегменте может объясняться низким качеством продукции, а в другом — высокими ценами. Следовательно, репозиторий моделей может включать специализированные модели для разных комбинаций цели, категории продукта и потребительского сегмента.

3.5.6. Таргетирование с использованием анализа выживаемости

Моделирование предрасположенности дает мощную основу для оценки вероятностей потенциальных результатов маркетинговых действий. Однако этот подход имеет ряд недостатков. Первая проблема заключается в том, что вероятность события не преобразуется напрямую во время-до-события, которое обычно является более эффективным показателем. Например, намного полезнее знать, что клиент, скорее всего, совершит покупку через 10 дней, и это время можно уменьшить на 5 дней, предложив скидку, чем знать, что условная вероятность покупки с учетом скидки равна 0,8. Обратите внимание, что эту проблему нельзя обойти, построив несколько моделей предрасположенности для смежных временных интервалов, потому что интервалы взаимозависимы. Например, нельзя построить отдельные модели вероятности покупки для января, февраля и марта, потому что покупки в феврале зависят от покупок в январе и т. д. Вторая проблема заключается в том, что мы не всегда имеем результаты, необходимые для создания меток откликов в моделировании предрасположенности. Например, можно обучить модель подобия для кампании удержания, чтобы различать клиентов, склонных и не склонных к оттоку. Набор обучающих данных будет включать профили клиентов, которые сбежали в течение некоторого периода в прошлом, а также профили клиентов, которые не сбежали в тот же период. Этот подход не идеален, потому что клиен-

ты, которые еще не сбежали, могут сбежать в будущем, поэтому точнее было бы говорить о том, что их результаты неизвестны, а не положительны. Это проблема цензурированных наблюдений, о которой мы уже говорили.

Эти ограничения моделирования предрасположенности можно преодолеть с помощью анализа выживаемости, представленного в разделе 2.6.2. Модели выживаемости способны правильно обрабатывать цензурированные данные, прогнозировать ожидаемое время-до-события (время выживаемости) и определять, как маркетинговые действия и свойства клиента способны ускорить или замедлить события. Рассмотрим численный пример, иллюстрирующий базовое использование анализа выживаемости в системе таргетирования.

ПРИМЕР 3.2

Рассмотрим сценарий, когда ретейлер настраивает рекламную кампанию в программной системе. Для определения оптимальных свойств кампании система использует набор данных, созданный для предыдущей аналогичной кампании, представленный в табл. 3.5. Он включает 12 профилей клиентов с 3 признаками: индикатор совершения покупки за неделю до объявления кампании, количество электронных писем, отправленных клиенту в рамках кампании, и размер скидки, предложенной клиенту. Наблюдаемый результат — это время покупки, измеряемое в днях, прошедших с момента объявления кампании. Кампания завершилась через 20 дней, поэтому все три клиента, которые не совершили покупки до окончания кампании, считаются цензурированными.

Мы используем этот небольшой набор данных, чтобы обучить модель Кокса пропорциональных рисков, как описано в разделе 2.6.2.3. Напомню, что модель Кокса является полупараметрической моделью с непараметрической базовой функцией выживаемости, которая описывает распределение времени покупок, и параметрической линейной моделью для оценки индивидуальных коэффициентов риска. Коэффициент риска описывает, будет ли «риск» покупки для данного клиента выше или ниже базового уровня. Коэффициент риска также выражается как функция признаков профиля, поэтому можно количественно определить, как разные признаки влияют на ожидаемое время покупки. Обучив модель Кокса, мы получаем следующую модель коэффициентов риска:

$$\begin{aligned}\log(risk) = & 1,957 \times \text{Предыдущая покупка} \\ & - 0,510 \times \text{Количество электронных писем} \\ & + 0,323 \times \text{Скидка}\end{aligned}\tag{3.20}$$

Таблица 3.5. Набор обучающих данных для анализа выживаемости.

Цензурированные записи соответствуют клиентам, не совершавшим покупок в течение первых 20 дней после начала кампании

№ п/п	Предыдущая покупка	Количество электронных писем	Скидка, %	Время покупки
1	0	2	5	5
2	0	2	0	10
3	0	3	0	20 (цензурировано)
4	1	1	0	6
5	1	2	10	2
6	1	3	0	15
7	1	4	0	20 (цензурировано)
8	1	5	5	6
9	0	2	10	8
10	1	5	5	13
11	0	0	0	20 (цензурировано)
12	1	2	5	8

Эту модель можно интерпретировать так: предыдущая покупка и скидка имеют отрицательную корреляцию с временем покупки, а количество электронных писем — положительную. Иначе говоря, дополнительная скидка сокращает время до покупки, а дополнительные электронные письма — увеличивают. Это означает, что отправка большего числа писем вредит кампании, поэтому стратегию использования электронной почты и необходимость отправки сообщений следует пересмотреть и исправить. Эта часть модели полезна, но пока не дает дополнительных сведений о стандартном моделировании предрасположенности. Большой интерес вызывают функции выживаемости. Модель Кокса может произвести функцию выживаемости для любого заданного вектора признаков, и каждая функция соответствует кумулятивному распределению времени покупок. В этом примере вектор признаков является трехэлементным вектором с индикатором предыдущей покупки, количеством электронных писем и величиной скидки в процентах. Примеры функций выживаемости представлены на рис. 3.14 и 3.15. Все кривые имеют одинаковую форму, но масштабированы в соответствии с коэффициентом риска, рассчитанным на основе вектора признаков x . Можно

видеть, что количество электронных писем подталкивает кривую выживаемости вверх, подтверждая и количественно оценивая неэффективность коммуникаций. Между тем скидка толкает кривую вниз, на уменьшение времени до покупки.

Важно не только получить кривые выживаемости, но и оценить статистические свойства времени до покупки. Напомню, что функции выживаемости $S(t)$ напрямую связаны с кумулятивными функциями распределения $F(t)$ времени до покупки:

$$S(t) = 1 - F(t). \quad (3.21)$$

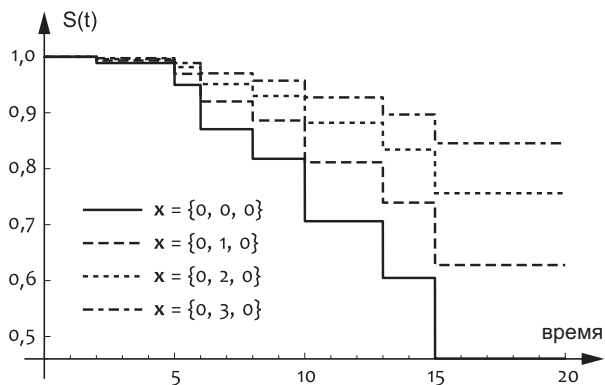


Рис. 3.14. Кривые выживаемости для разного количества электронных писем. Индикатор покупки и величина скидки равны нулю для всех кривых

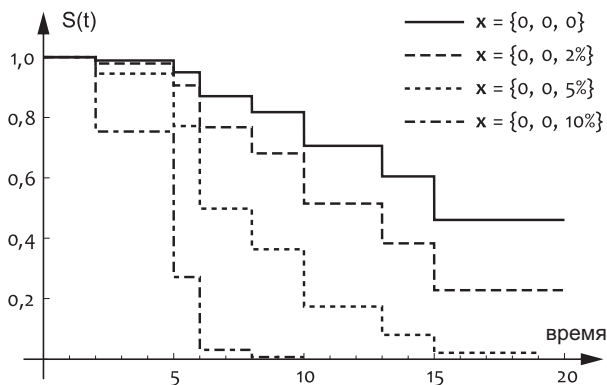


Рис. 3.15. Кривые выживаемости для разных скидок. Индикатор покупки и количество электронных писем равны нулю для всех кривых

Следовательно, из модели Кокса мы также получаем функции распределения. Это позволяет оценить среднее, или медианное, время до покупки, а также доверительные интервалы и другие статистические свойства. Поскольку функцию выживаемости можно получить для любой комбинации независимых переменных, мы получаем возможность оценить среднее, или медианное, время до покупки для каждого клиента отдельно, а затем использовать эти значения в маркетинговых правилах (например, отправить уведомление за день до предполагаемого времени покупки) или для нужд таргетирования (например, уделить внимание десяти процентам клиентов с наибольшим ожидаемым временем до покупки). Также можно количественно оценить влияние независимых переменных с точки зрения среднего или медианного времени до покупки. Например, можно оценить, как среднее количество дней до покупки уменьшается с каждым процентом скидки.

Анализ выживаемости можно применить ко многим маркетинговым мероприятиям. К наиболее типичным случаям можно отнести оценку правильности выбора времени для отправки сообщений в кампаниях привлечения, оценку времени до оттока в кампаниях удержания и оценку общего числа покупок в течение некоторого интервала времени для моделирования пожизненной ценности. Как и модели предрасположенности, модели выживаемости можно создавать для разных продуктов, категорий и клиентских сегментов. Ожидаемое время до покупки, полученное этими моделями, можно сравнить друг с другом, а затем выбрать наиболее релевантные продукты и предложения на основе соотношения времени.

3.5.7. Моделирование пожизненной ценности

Последний строительный блок, который мы рассмотрим, — оценка *пожизненной ценности клиента* (Customer Lifetime Value), сокращенно LTV, CLV или CLTV. Целью LTV-моделирования является оценка общей суммы, которую бренд может получить от данного клиента в течение всего срока их отношений. Точная структура анализа LTV в значительной степени зависит от бизнес-модели бренда, но можно создать базовые модели LTV и настраивать их с учетом специфических для бренда условий прибыли и убытков.

Анализ LTV — важный строительный блок в разработке кампаний и управлении маркетинг-миксом. Модели таргетирования способны помочь выбрать правильных клиентов, но анализ LTV позволяет количественно оценить ожидаемый результат таргетирования с точки зрения доходов и прибыли. Важность LTV также обусловлена тем, что из него можно получить другие основные показатели и пороговые значения, необходимые для принятия решения. Например, LTV является есте-

ственным верхним пределом расходов на привлечение клиента, а сумма LTV всех клиентов бренда, известная как *клиентский капитал* (customer equity), является основным показателем оценки стоимости бизнеса. Как и во многих других задачах маркетинговой аналитики и алгоритмического маркетинга, к моделированию LTV можно подойти с описательной, прогнозирующей и предписывающей точек зрения. Начнем с базового, описательного подхода, а затем перейдем к более продвинутым моделям.

3.5.7.1. Описательный анализ

Обычно LTV учитывает все доходы, полученные от клиента, и переменные затраты, связанные с отношениями с ним, и может дополнительно включать затраты на его привлечение. Одним из основных способов оценки пожизненной ценности клиента u является суммирование средней ожидаемой прибыли за некоторый промежуток времени в будущем:

$$\text{LTV}(u) = \sum_{t=1}^T (R - C) = T(R - C), \quad (3.22)$$

где время t измеряется в некоторых единицах, обычно месяцах, R и C — средние ожидаемые доходы и затраты соответственно на одного клиента в единицу времени, а T — ожидаемый срок жизни, или горизонт прогноза. Средние ожидаемые доходы и затраты обычно оцениваются по историческим данным, таким как история сделок и бюджеты кампаний. Эта оценка не является персонализированной (она усредняется по всем клиентам) и поэтому вычисляется относительно просто. Доходы и затраты могут резко различаться для разных сегментов клиентов, поэтому очень часто R и C оцениваются отдельно для каждого сегмента, а затем, если известен сегмент клиента, вычисляются значения LTV для сегментов. Продолжительность жизни T также можно выбирать эвристически, на основе типичной продолжительности отношений, или горизонта планирования, часто 24 или 36 месяцев.

Базовая формула LTV 3.22 не учитывает некоторых основных эффектов. Во-первых, она явно не учитывает расходы на удержание клиентов. Мы можем скорректировать временной горизонт T согласно средней продолжительности отношений с клиентом, однако было бы удобнее включить в формулу коэффициент удержания R как параметр. Например, годовой коэффициент удержания 0,8 означает, что 20 % нынешних клиентов прекратят отношения в течение года. Во-вторых, LTV обычно измеряется в течение относительно длительного промежутка времени, 2–3 лет, поэтому, возможно, придется учесть тот факт, что деньги в настоящем стоят больше, чем та же сумма в будущем, например, введением учетной ставки d . Учетная ставка отражает стоимость привязки капитала на определенный период времени. Например, годовая учетная ставка в размере 0,15 означает, что приведен-

ная стоимость в 1 рубль должна считаться эквивалентной 1,15 рубля, полученным в один год. Учет ставки дает *чистую приведенную стоимость* LTV. Принимая во внимание эти два фактора, чистую прибыль от клиента можно оценить как $(R - C)$ для первого периода времени, $(R - C) \cdot r / (1 + d)$ для второго и т. д., что в конечном итоге сводится к следующему определению LTV [Berger and Nasr, 1998]:

$$\text{LTV}(u) = \sum_{t=1}^T \frac{(R - C)r^{t-1}}{(1 + d)^{t-1}}. \quad (3.23)$$

Эта формула широко используется и может рассматриваться как стандартное определение LTV. Это выражение, конечно, не включает в себя все эффекты, которые можно найти в реальной жизни, и его можно расширить, чтобы отразить другие процессы и параметры, влияющие на LTV. Например, чистую прибыль m можно смоделировать не просто как постоянное значение $R - C$, а как значение, постепенно увеличивающееся с течением времени по мере развития отношений с потребителем:

$$m_t = m_0 + (m_M - m_0)(1 - e^{-kt}), \quad (3.24)$$

где m_0 — чистая прибыль в начале отношений, m_M — потенциальный максимум прибыли, а $k = \ln(2)/\tau$ — скорость роста прибыли, указанная в терминах времени на полпути к максимальному значению. Время половины пути τ определяет, насколько быстро прибыль приближается к потенциальному максимуму — для каждой единицы времени τ разница между текущим значением прибыли и максимумом уменьшается вдвое. Чистую прибыль m_t можно вставить в выражение 3.23 вместо постоянного значения $R - C$.

ПРИМЕР 3.3

Продолжим изучение рассмотрением численного примера вычисления LTV. Допустим, модель имеет следующие параметры:

- чистая прибыль с начала отношений — $m_0 = 100$,
- потенциальный максимум прибыли — $m_M = 150$,
- время половины пути к максимальному значению прибыли — $\tau = 3$ года.
- коэффициент удержания — $r = 0,9$.
- учетная ставка — $d = 0,1$.

Подставив эти параметры в уравнения 3.23 и 3.24, получим результат, представленный в табл. 3.6. Номинальная чистая прибыль в первом столбце рас-

тет в соответствии с уравнением 3.24 и проходит полпути (125) через три года. Ожидаемая чистая прибыль является произведением номинальной чистой прибыли на общий коэффициент удержания r^{t-1} . Наконец, дисконтированная чистая прибыль получается умножением ожидаемой чистой прибыли на дисконтный множитель $(1 + d)^{t-1}$, LTV — это сумма годовой дисконтированной чистой прибыли.

Обратите внимание, что этот анализ не только дает общую пожизненную ценность (LTV), но и показывает ее динамику с течением времени. По частичным суммам дисконтированной чистой прибыли за один, два и более лет можно нарисовать кривую LTV по времени. Если кривая быстро насыщается, это означает, что большая часть значения извлекается в начале отношений и длинные отношения не принесут много дополнительной прибыли. Если кривая неуклонно растет в течение длительного времени, это означает, что поступления от клиента остаются прибыльными в долгосрочной перспективе.

Таблица 3.6. Пример расчета LTV для пятилетнего горизонта

Год	Чистая прибыль	Коэффициент удержания	Ожидаемая чистая прибыль	Дисконтный множитель	Дисконтированная чистая прибыль
1	100,00	1,00	100,00	1,00	100,00
2	110,31	0,90	99,28	0,91	90,26
3	118,50	0,81	95,99	0,83	79,33
4	125,00	0,73	91,13	0,75	68,46
5	130,16	0,66	85,40	0,68	58,33
LTV					396,38

Описательная модель LTV напоминает анализ RFM, просто экстраполируя в будущее средние доходы, наблюдавшиеся в прошлом. Он допускает некоторый уровень персонализации, если рассчитывается для отдельных клиентских сегментов, но не предсказывает, как свойства клиента и маркетинговые действия могут повлиять на пожизненную ценность.

3.5.7.2. Цепи Маркова

Описательная модель LTV не дает большой гибкости, когда речь идет о сложных циклах взаимодействий с клиентами с несколькими состояниями привлечения,

максимизации и удержания. В то же время наличие нескольких состояний позволяет предположить возможность смоделировать цикл взаимодействия как случайный процесс или, точнее, цепь Маркова. Идея состоит в том, чтобы определить набор состояний клиента на основе наблюдаемых свойств, таких как давность покупки, оценить вероятность перехода между различными состояниями с соответствующими доходами и убытками, а затем оценить LTV на основе ожидаемого пути клиента в графе состояний [Pfeifer and Carraway, 2000].

Ключевой частью подхода на основе цепей Маркова является определение состояний и переходов, поэтому опишем его на примере. Допустим, ретейлер по своим данным определил, что давность последней покупки является хорошим индикатором оттока клиентов — клиенты, совершившие покупку в прошлом месяце, совершат покупку и в следующем с вероятностью $p_1 = 0,8$. Клиенты, совершившие последнюю покупку два месяца назад, совершат покупку еще раз с вероятностью $p_2 = 0,4$, давность в три месяца соответствует вероятности $p_3 = 0,1$, и, наконец, клиенты, не проявлявшие активности в течение четырех месяцев, едва ли вернутся вновь. Эту закономерность можно смоделировать с помощью цепи Маркова, как показано на рис. 3.16. Цепь имеет четыре состояния — по одному для каждого значения давности и одно для сбежавших клиентов. Клиенты, не совершавшие покупок, перемещаются по цепочке слева направо, шаг за шагом, пока не достигнут неактивного состояния. Покупка повторно инициализирует процесс и возвращает клиента в исходное состояние.

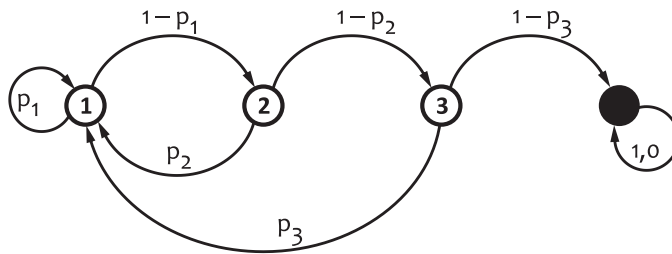


Рис. 3.16. Пример цепи Маркова для моделирования LTV. Белые круги соответствуют трем разным значениям давности. Черный круг представляет неактивное состояние

Эта цепь Маркова соответствует следующей матрице переходов:

$$P = \begin{bmatrix} p_1 & 1-p_1 & 0 & 0 \\ p_2 & 0 & 1-p_2 & 0 \\ p_3 & 0 & 0 & 1-p_3 \\ 0 & 0 & 0 & 1 \end{bmatrix}. \quad (3.25)$$

Каждая строка матрицы соответствует текущему состоянию, а каждый столбец — следующему. Элементы матрицы описывают вероятность перехода клиента из текущего состояния в следующее. Вычислить вероятность, что клиент, в настоящее время находящийся в состоянии s , через t месяцев окажется в состоянии q , можно как элемент (s, q) матрицы P^t , согласно стандартным свойствам цепи Маркова. Это дает простой способ оценить цикл взаимодействий с клиентом в вероятностных терминах, если известно текущее состояние.

С экономической точки зрения каждому состоянию цепи соответствуют прибыли и издержки. Например, маркетинговая стратегия может заключаться в том, чтобы потратить некоторый бюджет C на каждого активного клиента (например, послать печатный каталог) и прекратить делать это после перехода клиента в неактивное состояние. Первое состояние также связано с доходом R покупки. Введем вектор столбца \mathbf{G} , i -му элементу которого соответствует чистая прибыль i -го состояния:

$$\mathbf{G} = \begin{bmatrix} R - C \\ -C \\ -C \\ 0 \end{bmatrix}. \quad (3.26)$$

Очевидно, что матричное произведение $P\mathbf{G}$ даст в результате вектор ожидаемых чистых прибылей для каждого состояния через один период времени. Например, ожидаемая чистая прибыль для первого состояния будет:

$$\mathbb{E}[\text{profit} \mid \text{state } 1] = p_1(R - C) - (1 - p_1)C = p_1 \cdot R - C. \quad (3.27)$$

потратив сумму C , мы получаем шанс p_1 получить прибыль. Аналогично ожидаемая прибыль за второй период времени оценивается как $\mathbf{P}^2\mathbf{G}$ и т. д. Следовательно, LTV можно оценить как сумму таких ожидаемых прибылей за несколько периодов времени, а если дополнительно скорректировать каждый период по учетной ставке d , мы получим следующее выражение:

$$\mathbf{V} = \sum_{t=1}^T \frac{1}{(1+d)^{t-1}} \mathbf{P}^t \mathbf{G}, \quad (3.28)$$

где вектор столбца \mathbf{V} содержит оценки LTV для каждого начального состояния. LTV клиента также оценивается как один из элементов этого вектора на основе текущих состояний клиента, то есть в данном примере — значения давности. Этот результат можно сравнить со стандартной описательной моделью LTV в уравнении 3.23 — мы, по сути, заменили статические параметры с чистой прибылью и коэффициентом удержания вероятностной оценкой, зависящей от времени.

Завершим пример оценкой LTV для нескольких разных значений временного горизонта T . Как и прежде, возьмем за основу вероятности перехода $p_1 = 0,8$, $p_2 = 0,4$ и $p_3 = 0,1$. Предположим также, что ожидаемый доход от одной покупки составляет $R = 100$, ежемесячная стоимость маркетинговых коммуникаций $C = 5$, а месячная учетная ставка $d = 0,001$. Вычисляя выражение 3.28 для этих параметров и разных значений временного горизонта T , получаем следующую последовательность векторов LTV:

$$V_{T=1} = \begin{bmatrix} \$75,0 \\ \$35,0 \\ \$9,5 \\ \$0,0 \end{bmatrix} \quad V_{T=2} = \begin{bmatrix} \$135,5 \\ \$48,4 \\ \$10,4 \\ \$0,0 \end{bmatrix} \quad V_{T=3} = \begin{bmatrix} \$184,0 \\ \$53,4 \\ \$10,5 \\ \$0,0 \end{bmatrix}. \quad (3.29)$$

Как видите, LTV сильно зависит от исходного состояния клиента. Для временного горизонта в три месяца, то есть $V_{T=3}$, LTV клиента, сделавшего покупку месяц назад, составляет 184,0. Через два месяца ожидаемый LTV упадет до 53,4 и, наконец, после третьего месяца снизится до 10,5.

Метод на основе цепи Маркова можно расширить и приспособить для более сложных состояний клиента и маркетинговых стратегий. Например, вероятность следующей покупки часто коррелирует не только с давностью, но и с частотой покупок в прошлом. В этом случае каждую пару значений давности и частоты можно смоделировать как отдельное состояние в цепочке. Алгебраическую оценку LTV, заданную выражением 3.28, можно также заменить аппроксимацией цикла взаимодействий с клиентом методом Монте-Карло, что может дать еще большую гибкость в моделировании доходов и расходов. В этом случае мы случайно выбираем начальное состояние в соответствии с частотами, полученными по данным; затем пересекаем граф, подбрасывая монету в каждом состоянии, чтобы выбрать направление дальнейшего движения, и записываем доходы и расходы, возникающие на этом пути. Многократно повторяя этот процесс, получаем несколько вариантов ожидаемой пожизненной ценности. Преимущество этого подхода заключается в том, что статистические свойства LTV, такие как среднее значение, дисперсия и доверительные интервалы, можно напрямую оценить путем анализа гистограммы полученных результатов. Это также позволит включить в каждое состояние дополнительную бизнес-логику и параметры, что довольно сложно смоделировать с помощью матриц перехода.

3.5.7.3. Регрессионные модели

Цепи Маркова улучшают описательную модель LTV, заменяя статический коэффициент удержания и среднюю ожидаемую прибыль оценками, зависящими от времени и состояния. Однако этот подход имеет существенное ограничение —

число состояний растет экспоненциально с числом свойств клиента, включаемых в модель. Мы можем отступить на шаг назад и отметить, что концептуально как описательные модели, так и модели на основе цепей Маркова оценивают LTV с точки зрения вероятности, что клиент останется верным бренду и мы получим ожидаемую чистую прибыль. Это можно выразить следующим образом:

$$\text{LTV}(u) = \sum_{t=1}^T p(u, t) \cdot m(u, t), \quad (3.30)$$

где $p(u, t)$ — вероятность, что клиент u останется верным бренду до момента времени t , а $m(u, t)$ — чистая прибыль от клиента в период времени до момента t . Обычные описательные модели оценивают оба фактора с помощью статических значений коэффициента удержания и средней прибыли, тогда как модель на основе цепи Маркова оценивает те же факторы с применением вероятностного анализа. Более гибкое решение задачи 3.30 можно получить, создав регрессионные модели для обоих факторов. Преимущество этого подхода заключается в том, что регрессионные модели могут использовать широкий спектр независимых переменных, созданных на основе профиля клиента, и, таким образом, позволяют использовать предиктивные и предписывающие возможности.

Очевидно, что анализ выживаемости является естественным выбором для фактора вероятности удержания в уравнении 3.30. Эта вероятность напрямую соответствует функции выживаемости клиента $S_u(t)$, поэтому модель можно переписать как

$$\text{LTV}(u) = \sum_{t=1}^T S_u(t) \cdot m(u, t). \quad (3.31)$$

Модель выживаемости обучается для оценки времени оттока и требует определения события оттока. Эти события можно отслеживать непосредственно (если клиент явно отказывается от подписки на услугу) или эвристически, с помощью некоторых бизнес-правил (например, всех клиентов, не проявлявших активности пять и более месяцев, можно считать потерянными). Анализ выживаемости эффективно решает задачу оценки вероятности удержания путем правильной обработки цензурированных данных, а способность оценивать персонализирует функции выживания, параметризованные свойствами клиента, такими как давность и частота покупок.

Значения чистой прибыли $m(u, t)$ можно оценить несколькими разными способами с разной точностью. Одним из наиболее простых приближений является оценка среднего значения чистой прибыли для каждого сегмента клиентов и использование этого статического значения для всех клиентов в сегменте. Более сложные регрессионные модели можно создать включением параметра сезонности и признаков из профиля клиента.

3.6. Проектирование и проведение кампаний

Модели таргетирования и оценки LTV как основные строительные блоки системы таргетирования обеспечивают прочную основу для принятия эффективных маркетинговых решений. Однако маркетинговая кампания, как правило, представляет собой поток действий и решений, направленных на достижение определенной цели. Этот поток может потребовать связать несколько моделей и оптимизировать их с учетом множества сигналов и ограничений. Система таргетирования часто имеет некоторый репозиторий для шаблонов кампаний, где каждый шаблон описывает определенный поток действий и решений. Каждый поток обычно преследует определенную цель, но его можно параметризовать разными моделями таргетирования, размерами бюджета, свойствами пользователей и т. д. Прогнозирование рентабельности кампании и оптимизация параметров и пороговых значений для балансировки затрат и прибыли являются важной частью процесса проектирования кампании, а соответствующие процедуры и модели можно рассматривать как часть шаблона. В этом разделе мы рассмотрим несколько типов кампаний и их связь с базовыми моделями, о которых говорилось выше.

3.6.1. Цикл взаимодействий с клиентом

С экономической точки зрения, взаимодействие бренда с клиентом можно рассматривать как совокупность транзакций, характеризующихся общими суммами, приобретенными товарами, доходами, переходами по ссылкам на веб-сайте и т. д. Задачу оптимизации маркетинга тоже можно рассматривать с позиции транзакций, чтобы все компоненты маркетинг-микса были направлены на оптимизацию отдельных транзакций с точки зрения их вероятностей и доходности. Понятие жизненного цикла клиента расширяет контекст этой оптимизации, но по-прежнему фокусируется на задачах и целях бренда, а не на опыте клиента. Этот подход является неполным во многих маркетинговых средах, включая розничную торговлю, где взаимодействие с клиентом ориентировано на его опыт, а успех бренда определяется его способностью обеспечить долгосрочное поддержание превосходного опыта, а не оптимизировать отдельные транзакции.

Один из популярных подходов к анализу и моделированию нужд клиента основан на понятии *карт путешествий клиентов* (Customer Journey Map, CJM). Карта путешествий рассказывает историю взаимоотношений между клиентом и брендом. Она может описывать всю дугу взаимодействий аналогично кривой жизненного цикла или фокусироваться на определенной области, например на одной покупке. Карта обычно изображается в виде диаграммы с шагами или этапами взаимодействия и переходами между ними. На рис. 3.17 показан упрощенный пример карты

путешествия клиента. Она отражает поток одной транзакции, но помещает его в контекст клиентского опыта и долгосрочных взаимодействий с брендом.

Путешествие начинается с триггеров, к которым относятся поиски новых идей и продуктов, подготовку к особым мероприятиям, таким как дни рождения, получение рекламы по электронной почте или необходимость пополнения расходного продукта. За триггером следует изучение информации о продукте и выбор канала покупки. Затем взаимодействие продолжается в рамках выбранного канала, включая выбор определенного продукта и его покупка, и завершается действиями после покупки, такими как возврат продукта или написание отзыва. В реальной жизни карты путешествий клиентов обычно намного сложнее и включают множество подробностей о поведении клиентов и процессе принятия решений, распределении клиентов по разным состояниям и переходам и т. д. Кроме того, карты путешествий клиентов часто создаются отдельно для каждого сегмента, потому что маршруты путешествий в разных сегментах могут существенно различаться.

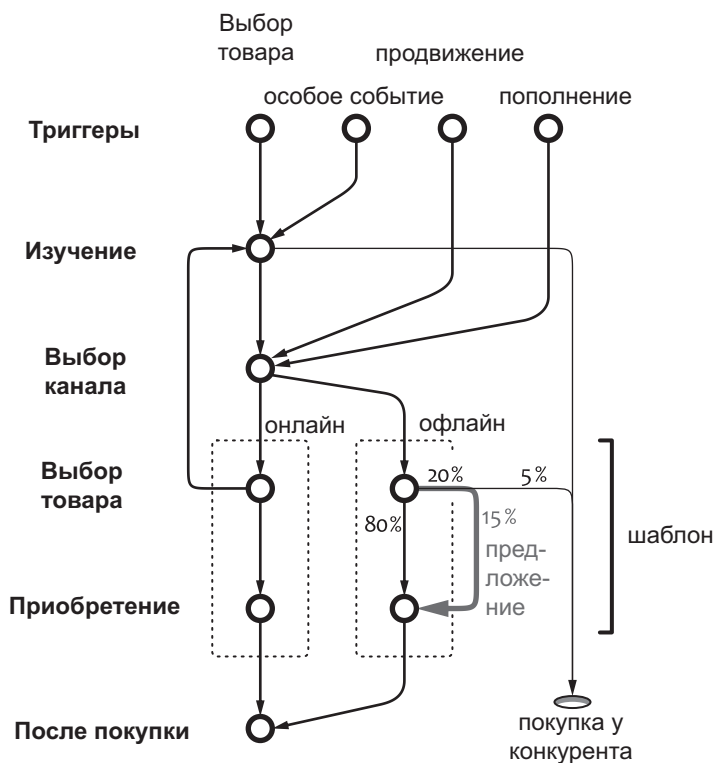


Рис. 3.17. Пример карты путешествия клиента

Маркетинговая кампания обычно имеет определенный след на карте путешествия, в том смысле что кампания пытается повлиять на путь, выбираемый клиентом. Например, на рис. 3.17 клиентам, предпочитающим производить покупки через интернет, делается выгодное предложение с целью вернуть их в обычные магазины. Каждую кампанию можно рассматривать как шаблон, который можно применить к конкретной ситуации в путешествии клиента. Программная система может иметь репозиторий *шаблонов кампаний*, где каждый шаблон включает правила, определяющие порядок инициирования действий в рамках кампании и реакции на ситуации, а также модели для оценки параметров необходимых действий и прогнозирования результатов. Шаблон может описывать одно или целый набор действий, которые могут выполняться в разные моменты времени, с использованием разных каналов и с учетом наблюдаемой обратной связи.

Анализ циклов взаимодействий с клиентами и создание карт путешествий обычно является стратегическим проектом, часто включающим обширные аналитические исследования, опросы клиентов и разработку маркетинговой стратегии. То есть создание карт путешествий клиентов и шаблонов кампаний не является обязанностью программной системы. Обычно предполагается, что эти артефакты создаются в другом месте и затем вводятся в систему. В обязанности системы, однако, входят оценка и оптимизация параметров шаблона и динамический выбор наиболее оптимальных шаблонов.

3.6.2. Кампании по продвижению продуктов

Один из основных видов целевых кампаний — стимулирование продаж отдельного товара. Примерами могут служить: реклама без указания денежной стоимости, купоны на скидку, купоны «два по цене одного» и бесплатные образцы продукции. В области торговли товарами широкого потребления кампании этого вида часто называют *рекламными вкладышами* (free-standing inserts, FSI) из-за буклетов купонов, которые вставляются в местные газеты. В самом простом виде отдельная рекламная акция соответствует простому путешествию клиента с триггером (реклама) и покупкой (удержание). Как будет показано ниже, этот подход не всегда является самым эффективным, но применим ко всем задачам:

- Для привлечения бренд может отправлять купоны «два по цене одного» или со скидкой активным покупателям, которые не покупают этот конкретный бренд.
- Для максимизации доходов бренд может проводить условные акции, такие как «купи 3 и получи скидку 1 рубль».
- Для удержания бренд может отправлять купоны «два по цене одного» или со скидкой покупателям, уменьшившим потребление по сравнению с предыдущими циклами покупки.

Платформа моделирования ответов содержит некоторые рекомендации, как нацелить такие акции с помощью прогнозирующих моделей, которые оценивают вероятность отклика, но есть много дополнительных аспектов, которые необходимо предусмотреть, включая процесс таргетирования, правила составления бюджета и выбор параметров продвижения.

3.6.2.1. Процесс таргетирования

Систему таргетирования можно использовать как в пакетном режиме, так и в режиме реального времени, в зависимости от среды и характера кампании. Некоторые рекламные акции можно проводить, отправляя миллионы электронных писем за раз, поэтому система таргетирования может заранее подготовить список клиентов. Другие кампании проводятся практически в режиме реального времени из-за быстро меняющихся профилей клиентов или контекста. Например, рекламные предложения могут делаться клиентам в зависимости от содержимого их корзин, непосредственно перед оплатой на кассе. Поход на основе реального времени, как правило, получается более гибким, и правильно разработанная система таргетирования в реальном времени может имитировать также пакетный режим, оценивая правила и модели для всей базы данных клиентов. Поэтому сфокусируемся на случае таргетирования в реальном времени и рассмотрим процесс, который получает единственный профиль клиента и соответствующий контекст и создает список рекламных предложений для этого клиента.

Предположим также, что система имеет базу данных рекламных акций, которые могут быть предложены. Она включает рекламные акции из всех кампаний, активных в настоящее время. Каждая акция должна быть связана с такими свойствами, как бизнес-цель, продвигаемый продукт и категория, чтобы система таргетирования могла использовать эту информацию для выбора правильных моделей и правил таргетирования. Именно поэтому создание и таргетирование рекламных акций тесно связаны друг с другом, в том смысле что для каждого шага или признака в модели таргетирования должен иметься аналог в конфигурации кампании и атрибутах рекламных акций. Далее мы рассмотрим процесс таргетирования и обсудим, как выбираются рекламные акции из набора доступных вариантов, а также методологию создания рекламных акций и наделение их свойствами и условиями, необходимыми для таргетирования.

Процесс таргетирования можно рассматривать как последовательность из трех шагов. Сначала система извлекает все доступные акции и выбирает подходящие для данного контекста и клиента. Затем акции оцениваются и сортируются по степени пригодности для поставленной цели. Наконец, в соответствии с лимитом бюджета и другими ограничениями выбирается оптимальный набор, который будет предложен клиенту. Этот процесс показан на рис. 3.18. Начальная фильтра-

ция рекламных акций обычно основана на бизнес-правилах и условиях, поэтому мы называем ее *жестким таргетированием*. С другой стороны, для оценки акций обычно используются предиктивные модели, дающие непрерывную оценку, поэтому мы называем этот этап *мягким таргетированием*.

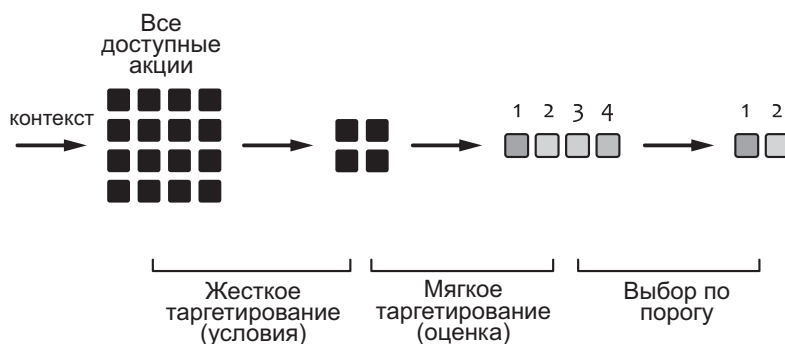


Рис. 3.18. Процесс таргетирования рекламных акций

Цель этапа жесткого таргетирования — выбрать рекламные акции, соответствующие данному контексту. Акции, создаваемые в системе таргетирования, обычно связаны с условиями, которым должен соответствовать данный контекст для активации акции. Целью условий является стимулирование определенного потребительского поведения и достижение основных экономических целей акции. Условия жесткого таргетирования, по существу, определяют шаблон кампании, то есть точку на пути клиента, где должна применяться акция. Рассмотрим следующие типичные примеры:

- Количественное условие. Активирует акцию, когда клиент покупает определенное количество определенного продукта, бренд или категорию за один переход или в течение определенного периода. Это условие часто используется в кампаниях максимизации, чтобы *подтянуть* потребителя, то есть дать стимул покупать больше обычного. Например, клиенту, который обычно покупает две пачки йогурта, можно предложить *купить 4 и получить 1 бесплатно*.
- Условие отсутствия покупок. Активирует акцию для клиентов, не покупавших товар или бренд в течение определенного периода времени. Это условие можно использовать в кампаниях по удержанию и привлечению, чтобы отделить активных потребителей бренда от неактивных и перспективных.
- Условие по типу канала. Активирует акцию, когда клиент взаимодействует с брендом или ретейлером по определенному каналу. Например, клиента можно вознаградить за три посещения магазина в неделю.

- **Условие перенацеливания.** Активирует акцию на основе ранее предложенных или использованных акций. Например, с клиентами, получившими, но не активировавшими рекламные акции, предложенные по цифровым каналам, можно связаться по каналам в магазине.
- **Условие по местоположению.** Активирует акцию, исходя из местоположения клиента, определяемого по данным мобильного устройства, местоположению магазина, маякам в магазине или IP-адресу.
- **Условие доступности.** Некоторые акции можно временно деактивировать, если соответствующие товары отсутствуют на складе или недоступны через данный канал сбыта.

Этап жесткого таргетирования выбирает рекламные акции, которые потенциально можно предложить потребителю. Цель этапа мягкого таргетирования — выбрать наиболее релевантные предложения и отфильтровать варианты, которые могут оказаться неэффективными. Мягкое таргетирование часто осуществляется с помощью моделей предрасположенности. Система таргетирования может иметь репозиторий моделей, где каждая модель обучается для определенной бизнес-цели и категории продукта и характеризуется соответствующими признаками. Поскольку акции обладают похожими свойствами, система может динамически связывать модели с акциями. Оценочные модели можно комбинировать со специальными условиями, дополняющими логику модели. Например, базовая модель привлечения, реализующая метод подобия, идентифицирует клиентов, похожих на естественных любителей пробовать новое, но не гарантирует, что рекламное предложение не будет сделано тем, кто и так покупает продукт. Напротив, предложения, направленные на максимизацию и удержание, как правило, не должны делаться клиентам, которые не потребляют рекламируемый продукт. Такие дополнительные проверки можно реализовать в виде условия.

3.6.2.2. Бюджет и ограничения

После подготовки и ранжирования кандидатов система должна выбрать окончательный набор акций, которые можно предложить клиенту. Этот шаг может включать определение разных ограничений, применяемых к кампании. Во-первых, следует ограничить количество рекламных акций, предлагаемых клиенту в рамках одной кампании, а также общую частоту общения с клиентом (количество сообщений в единицу времени). Эти правила, часто называемые *правилами давления* или *правилами ограничения частоты*, обычно используют пороги, выбранные эвристически или экспериментально. Затем определяется ограничение на бюджет кампании или общее количество рекламных акций. Однако часто система таргетирования должна сама определить оптимальное количество рекламных акций, чтобы максимизировать рентабельность кампании. Может случиться так, что это число

окажется намного ниже предела, указанного маркетологом, и полное расходование бюджета может привести к убыткам. С точки зрения моделирования предрасположенности проблему оптимизации прибыльности можно рассматривать как поиск порога предрасположенности, который максимизирует прибыль, когда все клиенты с оценкой выше порога рассматриваются как целевые, а все другие — нет. Мы уже видели, как с помощью моделирования ответов можно найти компромисс между затратами на кампанию и прибылью, а теперь рассмотрим пример, иллюстрирующий практические детали.

ПРИМЕР 3.4

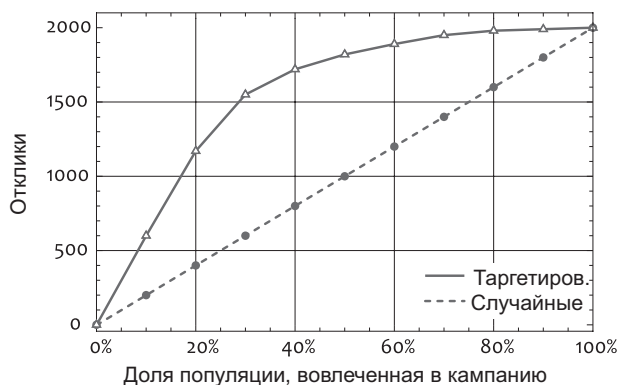
Рассмотрим случай с ретейлером, имеющим 100 000 держателей карт постоянного клиента. Ретейлер планирует целевую кампанию, каждая акция которой стоит 1 рубль, а потенциальная прибыль от одного отклика составляет 40 рублей. Средняя частота откликов для данного типа кампании и категории продуктов, оцененная по историческим данным, составляет 2 %. Имея модель предрасположенности, оценивающую вероятность отклика для каждого клиента, мы можем оценить и отсортировать всех держателей карт. Обобщая результаты, можно разделить клиентов на группы одинакового размера, где первая группа включает клиентов с наивысшими оценками, а последняя — с наименьшими. В этом случае задачу таргетирования можно определить как поиск оптимального количества групп с высшими оценками для включения в список таргетирования, что эквивалентно поиску пороговой оценки, отделяющей верхние группы от нижних. Группировка в этом примере используется только для удобства, и, хотя этот подход часто используется на практике, ничто не мешает выполнить те же вычисления для отдельных клиентов, то есть определить число групп по числу клиентов. Предположим, что у нас есть 10 групп, или децилей, то есть каждая группа включает 10 000 клиентов; следовательно, среднее ожидаемое число откликнувшихся — 200 человек на группу. Другими словами, случайно разбросав клиентов по группам, мы с высокой долей вероятности получим по 200 откликов от каждой группы. Это число показано во втором столбце в табл. 3.7, а в третьем столбце содержится общее число откликнувшихся, которое в нижней строке достигает 2000, или 2 % от клиентской базы.

Далее предположим, что модель предрасположенности сгенерировала для каждой группы наименьшие оценки вероятности отклика, представленные в четвертом столбце. Умножив размер группы на эту вероятность, мы получим ожидаемое число откликов в случае целевого распределения, представленного в следующих двух столбцах. Разумеется, общее число откликнувшихся все так же составляет 2000 человек. Отношение между

Таблица 3.7. Пример определения прироста эффективности кампании

Де- цель	Откликов, случай- ное распределение		Вероят- ность	Откликов, таргетиро- ванное распределение		Подъем
	На группу	Всего		На группу	Всего	
1	200	200	0,060	600	600	3,00
2	200	400	0,057	570	1170	2,93
3	200	600	0,038	380	1550	2,58
4	200	800	0,017	170	1720	2,15
5	200	1000	0,010	100	1820	1,82
6	200	1200	0,007	70	1890	1,58
7	200	1400	0,006	60	1950	1,39
8	200	1600	0,003	30	1980	1,24
9	200	1800	0,001	10	1990	1,11
10	200	2000	0,001	10	2000	1,00

числом откликнувшихся в целевом и случайном распределениях называется *подъемом*, и это ключевой показатель, описывающий качество модели таргетирования. Подъем обычно визуализируется с помощью *диаграммы подъема*, как показано на рис. 3.19. На этой диаграмме изображены две линии, соответствующие совокупному числу ответов: прямая соответствует случайному распределению, а изогнутая — целевому.

**Рис. 3.19.** Диаграмма подъема для модели таргетирования

Чтобы определить количество целевых групп, необходимо оценить рентабельность кампании. Каждая акция стоит 1 рубль, поэтому стратегия случайного распределения невыгодна, потому что каждая группа вызывает потерю:

$$\begin{aligned}
 &40 \text{ рублей за отклик} \times 10\,000 \text{ клиентов} \times 0,02 \text{ частота отклика} - \\
 &- 1 \text{ рубль на клиента} \times 10\,000 \text{ клиентов} = \\
 &= -2000 \text{ рублей}
 \end{aligned}$$

Целевая кампания, напротив, будет прибыльной для первых трех групп из-за высокой доли откликнувшихся, как показано в табл. 3.8. Как видите, с увеличением числа групп, вовлеченных в кампанию, рентабельность этой кампании сначала увеличивается, а затем уменьшается и в конечном итоге становится отрицательной.

Рентабельность кампании достигает максимума при охвате трех групп, то есть 30 % населения. Это соответствует всем клиентам с оценкой предрасположенности выше 0,038. График изменения рентабельности целевой кампании показан на рис. 3.20. Обратите внимание, что максимально возможный бюджет, соответствующий рекламному предложению каждому клиенту, не максимизирует рентабельности. Напротив, это приводит к потере 20 000 рублей.

Таблица 3.8. Пример вычисления прибыльности кампании

Де- цель	Стои- мость	Выгода		Таргетированная рентабельность
		случайное распределение	таргетированное распределение	
1	10000	-2000	14000	14000
2	10000	-2000	12800	26800
3	10000	-2000	5200	32000
4	10000	-2000	-3200	28800
5	10000	-2000	-6000	22800
6	10000	-2000	-7200	15600
7	10000	-2000	-7600	8000
8	10000	-2000	-8800	-800
9	10000	-2000	-9600	-10400
10	10000	-2000	-9600	-20000

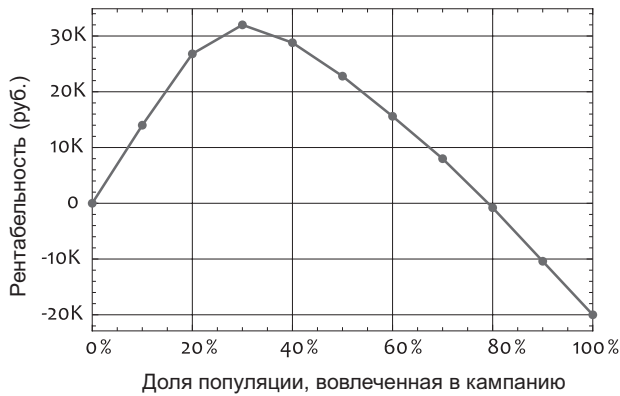


Рис. 3.20. Рентабельность таргетированной кампании как функция охвата

Важно отметить, что в этом примере мы использовали базовую вероятность отклика, а не моделирование увеличения эффективности. На практике это может привести к снижению эффективности кампании, поскольку высокий уровень откликов не гарантирует увеличение расходов или потребления клиентов. Другими словами, контрольная выборка в каждой группе может показывать одинаково хорошую отдачу или даже лучше, чем целевая выборка в той же группе. Эту проблему можно обойти, заменив вероятности отклика в табл. 3.7 оценками подъема, описанными в разделе 3.5.4.2.

Принцип максимизации рентабельности позволяет определить оптимальные базовые параметры кампании, такие как общее количество акций и порог рентабельности. В реальном мире иногда бывает полезно отклониться от базового уровня, особенно для приложений реального времени, когда набор клиентов, которые фактически будут взаимодействовать с системой, заранее неизвестен. Рассмотрим следующий сценарий. Система запускает рекламную кампанию с фиксированным бюджетом и равномерно распределяет этот бюджет в пределах временных рамок кампании. Это предполагает, что мы должны использовать какой-то фиксированный коэффициент распределения, например 100 промоакций в час. Однако что делать, если кампания уже достигла этого показателя (в данном примере уже было сделано 100 предложений в течение последнего часа) и мы сталкиваемся с потребителем с очень высоким показателем предрасположенности? В этом случае, может быть, разумно превысить бюджет, а затем немного уменьшить ставку, чтобы вернуться в прежнее русло. Такое поведение можно реализовать путем динамической корректировки пороговых значений при отклонении от целевой скорости распределения. Эта идея проиллюстрирована на рис. 3.21. Мы опреде-

ляем целевую скорость распределения и два значения — ε^- и ε^+ , — определяющие максимальное отклонение от целевой линии. Обратите внимание, что целевая линия не обязательно должна быть прямой, можно использовать более сложную кривую, учитывающую выходные дни, рабочие часы и т. д. Фактическая скорость распространения постоянно измеряется и контролируется системой, чтобы оставаться в пределах отклонений.

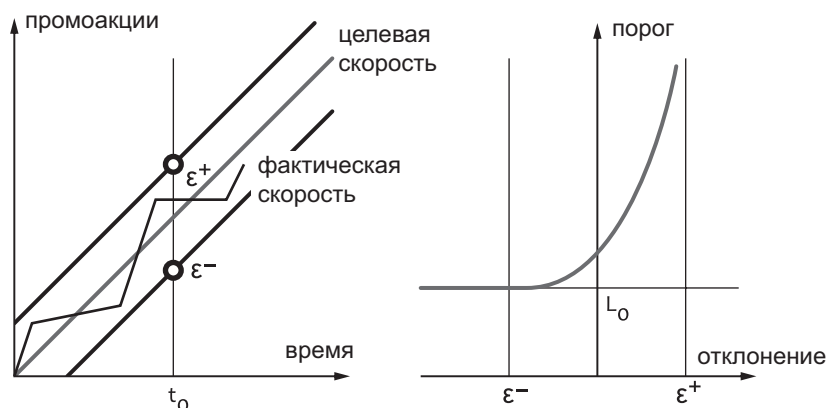


Рис. 3.21. Динамическая оценка порога для управления бюджетом

Оценку порога можно выразить как функцию отклонения от целевой линии в текущий момент времени t_0 . Если мы существенно ниже уровня бюджета (под линией ε^-), оценку порога можно уменьшить до минимума, что соответствует наибольшей близости L_0 потребителя и кампании, чего достаточно, чтобы сделать предложение. Если мы существенно превысили бюджет (над линией ε^+), оценку порога можно увеличить до максимума, чтобы полностью остановить распределение. Эти две крайние точки могут быть связаны некоторой растущей функцией, как показано на рис. 3.21. Как результат, по мере приближения и пересечения бюджетных лимитов мы становимся все более требовательными к потребителям и понижаем планку, когда перестаем сталкиваться с достаточно перспективными клиентами.

3.6.3. Многоступенчатые рекламные кампании

Отдельные рекламные кампании, такие как распространение купонов для привлечения и максимизации, широко используются на практике. Но такая стратегия не всегда эффективна, так как имеет очень короткое и ограниченное влияние на цикл взаимодействия с клиентом [Catalina Marketing, 2014]. Иногда можно создавать более сложные кампании, включающие несколько этапов и влияющие на поведе-

ние клиента в течение более длительного периода времени. Рассмотрим пример кампании, направленной на максимизацию продаж товаров широкого потребления, которая имеет следующую организацию:

- Первый этап кампании — ее объявление с целью информирования клиентов о предложении. Например, бренд может распространить через доступные маркетинговые каналы следующее сообщение: *Купите Q или более единиц продукта X и сэкономьте при следующем походе в магазин. Чем больше купите, тем больше сэкономите.*
- Второй этап — распределение. Система таргетирования отслеживает покупки и выдает купоны на скидку клиентам, удовлетворяющим условиям таргетирования, то есть в данном примере купившим Q и более единиц товара X . Величина скидки купона определяется динамически, в зависимости от количества купленного товара — чем больше потребитель купит, тем больше сэкономит. На этом этапе потребитель получает стимул купить больше единиц, чтобы получить купон в качестве вознаграждения.
- Третий и последний этап — погашение. Во втором походе по магазинам потребитель покупает рекламируемый продукт, чтобы погасить купон, полученный на предыдущем этапе. Потребитель получает стимул купить товар, погасить купон и получить скидку.

Этот шаблон кампании можно рассматривать как цикл взаимодействия с клиентом с тремя этапами: инициация, покупка и погашение. Можно утверждать, что этот подход более эффективен, чем отдельные рекламные акции, поскольку оказывает более длительное влияние на лояльность клиентов и более низкие затраты на единицу проданной продукции [Catalina Marketing, 2014]. Динамически определяемая величина скидки на втором этапе является интересной деталью — система таргетирования должна оптимизировать ее и предсказать, как она повлияет на результаты кампании. Этот аспект не учитывается в процессах определения целевых показателей и бюджета, о которых говорилось в предыдущих разделах. Рассмотрим пример, демонстрирующий, как система таргетирования может эвристически оценивать разные параметры продвижения и прогнозировать результаты кампании, используя только базовые статистики. Более формальные методы оптимизации скидок будут рассмотрены в главе 6, в контексте оптимизации цен.

ПРИМЕР 3.5

Рассмотрим случай рекламной кампании, которая следует сценарию с тремя этапами, описанному выше. Цель системы таргетирования — выбрать обоснованное пороговое значение количества Q и на втором этапе определить

размер скидки, исходя из фактически приобретенного количества. Начнем с гистограммы объемов покупок для продвигаемого продукта, рассчитанной для временного интервала, равного длительности кампании. Обозначим количество покупок ровно q единицами рекламируемого продукта как $H(q)$. Соответствующая историческая гистограмма выглядит следующим образом:

$$\begin{aligned} H(1) &= 4000 (32 \%) & H(4) &= 1000 (8 \%), \\ H(2) &= 5000 (40 \%) & H(5) &= 600 (5 \%), \\ H(3) &= 2000 (16 \%) & H(6) &= 0. \end{aligned}$$

Нам нужно подтянуть клиентов, покупающих продукт в относительно небольших количествах, поэтому система может выбрать порог Q выше большинства покупок. В данном примере разумным выглядит выбор числа 3, потому что 72 % покупок находится под этим порогом. Следовательно, система предложит купон на скидку клиентам, купившим больше 3 единиц. Величина скидки зависит от фактически приобретенного количества. Более подробно этот аспект мы обсудим в главе 6, а пока просто предположим, что величина скидки является статическим параметром конфигурации. Например, пусть минимальная скидка составляет 15 % и с каждым уровнем увеличивается на 5 %. То есть клиент, купивший 3 единицы, получит скидку 15 %, купивший 4 единицы получит 20 %, а купивший 5 единиц получит 25 %. Обозначим количество единиц на уровне i как q_i и соответствующее значение скидки как d_i . После определения всех этих параметров система может спрогнозировать результаты кампании. Это можно сделать отдельно для каждого уровня скидки. Ожидаемое количество купонов, сгенерированных на уровне i , можно определить на основе ранее созданной гистограммы как

$$coupons(i) = H(q_i). \quad (3.32)$$

Ожидаемое число погашений можно оценить с помощью модели откликов, использующей величину скидки как признак:

$$redemptions(i) = coupons(i) \times r(d_i), \quad (3.33)$$

где $r(d_i)$ — средняя доля откликов, предсказанная моделью. Стоимость купонов на уровне i можно оценить как

$$cost(i) = (\text{цена продукта} \times d_i \times q_i + c) \times redemptions(i), \quad (3.34)$$

где c обозначает дополнительные расходы, связанные с купоном, такие как расходы на распространение и услуги клиринговой организации. Эффектив-

ность кампании можно предсказать как соотношение между общим количеством погашений и суммарными затратами по всем уровням (стоимость за погашение).

3.6.4. Кампании по удержанию

Кампании по удержанию нацелены на удержание клиентов, которые могут уйти. Кампании этого типа широко используются в телекоммуникационных, страховых, банковских и других сферах, основанных на подписке, где непрерывность отношений имеет решающее значение. Однако проблема оттока клиентов актуальна и для других сфер, где нет подписки, включая розничную торговлю. Одна из ключевых причин, почему удержание так важно, заключается в том, что привлечь новых клиентов порой сложнее и дороже, чем удержать существующих. Согласно некоторым исследованиям, стоимость привлечения одного потребителя может быть в 10–20 раз выше стоимости удержания из-за более низкой частоты отклика и других факторов [Artun and Levin, 2015].

Кампанию по удержанию можно определить как последующую деятельность с клиентами, находящимися в группе риска. Однако сами действия и риски могут определяться по-разному, в зависимости от кампании. Примерами риска могут служить риск отмены подписки и риск миграции в другую сеть супермаркетов. Определение риска зависит от бизнес-модели, характера продукта или услуги и особенностей использования. Поставщика программных услуг, например, может беспокоить риск отмены подписки, а также значительное число клиентов, которые создают учетную запись, но не загружают клиентское приложение. Подобный риск незагрузки приложения можно распознать и устранить с помощью специальной кампании по удержанию. В числе примеров действий, направленных на это, можно назвать рассылку напоминаний по электронной почте, распространение учебных материалов, предложения оставить отзыв о недавно приобретенном продукте, а также специальные предложения и скидки.

В отличие от кампаний по продвижению продукта, при разработке кампаний по удержанию больше внимания уделяется пожизненной ценности и подъему. Включение прогнозируемой пожизненной ценности имеет большое значение, потому что бессмысленно инвестировать в удержание клиентов с низкой ценностью. Моделирование подъема важно, потому что ориентация на неправильных клиентов контрпродуктивна по нескольким причинам [Radcliffe and Simpson, 2007]. Во-первых, многие клиенты в группе риска уже недовольны, и дополнительные коммуникации, особенно навязчивые, такие как телефонные звонки, могут ускорить процесс оттока. Во-вторых, некоторые удерживающие коммуникации могут

напоминать клиентам, что у них есть возможность уйти, что заставит их пересмотреть свои отношения с брендом и искать альтернативные варианты. Поэтому важно, чтобы коммуникации были целенаправленными, а результаты постоянно оценивались с помощью контрольных групп.

Кампании удержания обычно собираются из стандартных блоков, но есть другие методы проектирования. Одним из самых основных подходов является таргетирование предрасположенности к оттоку. Эту модель можно создать, используя стандартные методы моделирования предрасположенности с обучающим набором данных, собранным из активных и утраченных профилей клиентов. Этот, казалось бы, простой подход имеет подводные камни, которые следует обсудить. Как уже говорилось выше, маркетинговое действие можно описать в терминах двух условных вероятностей — вероятности отклика при воздействии и вероятности отклика без воздействия. В случае с удержанием событию отклика соответствует отток. Всех клиентов можно классифицировать по этим двум вероятностям, как показано на рис. 3.22.

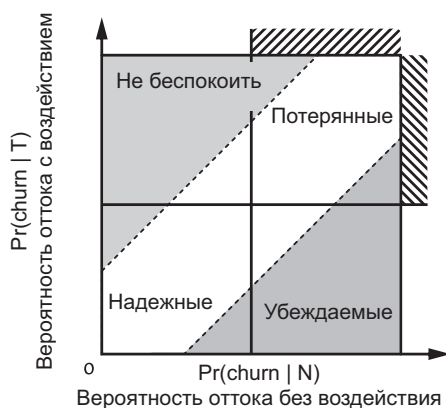


Рис. 3.22. Классификация клиентов с точки зрения кампании по удержанию

Если общая стратегия удержания сфокусирована, то есть обычно на клиентов не оказывается воздействий, направленных на удержание, модель предрасположенности, обученная предсказывать результаты оттока, фактически предсказывает вероятность оттока при отсутствии воздействия как

$$\text{score}(\mathbf{x}) = \text{Pr}(\text{churn} | N, \mathbf{x}), \quad (3.35)$$

где \mathbf{x} — вектор признаков профиля клиента. Кампания удержания, управляемая этой вероятностью, фокусируется на самом правом вертикальном срезе на рис. 3.22,

включающем не только много убеждаемых, но и много потерянных клиентов. Если стратегия удержания достаточно широка, то есть если почти все клиенты подвергаются некоторому воздействию, модель будет фактически оценивать предрасположенность к оттоку при наличии воздействия как

$$\text{score}(\mathbf{x}) = \Pr(\text{churn} | T, \mathbf{x}), \quad (3.36)$$

что соответствует горизонтальной области в верхней части квадрата на рис. 3.22, содержащей много потерянных клиентов и клиентов, не терпящих, когда их беспокоят. Этот аспект следует учитывать при выборе совокупности для обучения модели. Кампания по удержанию может также использовать анализ выживаемости для оценки времени оттока, который лучше подходит для выбора правильного момента воздействия, чем оценка вероятности оттока.

Таргетирование на основе вероятности оттока клиентов не учитывает долгосрочные результаты кампании. Эти результаты можно количественно оценить в терминах пожизненной ценности, потому что каждое удержание сохраняет LTV соответствующего клиента, а отток каждого клиента является потерей его LTV. Произведение вероятности оттока и оценки LTV для данного клиента определяет *ожидаемый убыток*. Если предположить, что воздействие на клиентов с наибольшими ожидаемыми потерями максимизирует отношение между сохраненными доходами и затратами на кампанию, тогда этот показатель можно использовать для таргетирования:

$$\text{score}(\mathbf{x}) = \Pr(\text{churn} | N, \mathbf{x}) \times \text{LTV}(\mathbf{x}). \quad (3.37)$$

Величину LTV можно оценить по средним расходам клиента или с помощью более продвинутых моделей LTV, описанных ранее. Эту модель можно настроить в зависимости от бизнес-модели для учета затрат и прибылей, связанных с различными возможными результатами. Например, можно отдельно оценить ожидаемые доходы от удержания, потери из-за оттока и затраты на кампанию. Модель ожидаемых потерь из уравнения 3.37 широко используется на практике из-за своей простоты и достаточно высокой эффективности.

Основным недостатком модели ожидаемых потерь является использование ею только вероятности, а не подъема оттока, то есть разницы между вероятностями оттока клиентов, подвергавшихся и не подвергавшихся воздействию:

$$\text{uplift}(\mathbf{x}) = \Pr(\text{churn} | T, \mathbf{x}) - \Pr(\text{churn} | N, \mathbf{x}). \quad (3.38)$$

Положительный подъем оттока означает, что воздействие усиливает отток, то есть воздействие оказывает отрицательный эффект. Высокий подъем соответствует

верхнему левому углу квадрата на рис. 3.22. Отрицательный подъем оттока означает, что воздействие уменьшает отток, что соответствует правому нижнему углу на рис. 3.22. Следовательно, мы должны ориентироваться на клиентов, используя оценку, обратную подъему:

$$\text{score}(\mathbf{x}) = -\text{uplift}(\mathbf{x}) = \text{savability}(\mathbf{x}). \quad (3.39)$$

Эту метрику называют также оценкой *сохраняемости* (savability), потому что она оценивает склонность положительно реагировать на действия по удержанию. Подъем/сохраняемость можно смоделировать с помощью методов, описанных в разделе 3.5.4.2, включая подходы с одной или двумя моделями. Так же как в других приложениях моделирования подъема, подход на основе оценки возможности спасения помогает отделить клиентов, которые наверняка останутся, только если воздействовать на них, и тем самым повысить эффективность кампании удержания. Но имейте в виду, что этот подход также наследует типичные недостатки моделирования подъема. К ним относятся повышенная сложность моделирования и более высокая дисперсия оценок, потому что подъем — это разность двух случайных величин. Подъем также можно объединить с методом оценки ожидаемых потерь, чтобы учесть долгосрочное воздействие на полученные LTV:

$$\text{score}(\mathbf{x}) = \text{savability}(\mathbf{x}) \times \text{LTV}(\mathbf{x}). \quad (3.40)$$

После вычисления оценок оптимальную глубину таргетирования, то есть процентную долю целевой популяции, можно определить с помощью метода максимизации рентабельности, описанного выше, в разделе 3.6.2.2. Затем кампанию можно провести с тем же процессом таргетирования, что и в кампаниях по продвижению продукта.

3.6.5. Кампании пополнения

Кампании по удержанию наиболее актуальны для компаний, основанных на подписке, таких как телекоммуникационные услуги, страхование, ПО и банковское дело. В сфере розничной торговли модель подписки используется реже, но многие продукты приобретаются на регулярной основе, поэтому модель взаимодействия становится похожей на подписку. Примеров пополняемых продуктов можно привести массу, это и продукты питания, и косметика, и канцелярские принадлежности, и такие аксессуары, как фильтры для воды, и многое, многое другое. Кампании пополнения направлены на стимулирование повторных покупок и уменьшение *цикла покупок* отправкой напоминаний, рекомендаций и специализированных предложений.

С точки зрения разработки, отличительными особенностями кампаний пополнения является акцент на времени коммуникации и покупательских привычках. Время коммуникации играет важную роль, потому что уведомления о пополнении должны согласовываться с циклами покупок — например, бессмысленно посылать уведомление клиенту сразу после приобретения продукта. Связь с покупательскими привычками также важна, поскольку текст уведомления должен соответствовать продуктам и категориям, обычно приобретаемым его получателем.

Начнем с очень простого подхода, который может быть реализован системой таргетирования. Сначала система оценивает среднюю продолжительность цикла закупок для каждого пополняемого продукта или категории продуктов. Затем снова и снова, например ежедневно, запускается кампания пополнения. Каждый раз система просматривает профили клиентов и определяет последнюю дату покупки пополняемых продуктов. Эта дата сравнивается с предполагаемой продолжительностью цикла закупок, и клиентам, близким к концу цикла, посылается уведомление. Сообщение может быть персонализировано на основе последних или наиболее часто приобретаемых продуктов, найденных в истории покупок.

Одним из основных ограничений этого подхода является слишком грубая оценка циклов пополнения. Ее можно улучшить, если, например, разбить оценки не только по категориям продуктов, но также по сегментам клиентов, чтобы учесть различия между клиентами. Другими словами, продолжительность цикла оценивается для каждой пары категория/сегмент. Еще более точные результаты можно получить, используя для оценки времени покупки анализ выживаемости. Модель выживаемости также позволяет определить факторы, положительно или отрицательно влияющие на время покупки, такие как скидки или уведомления о пополнении, поэтому содержание и частоту сообщений можно соответствующим образом скорректировать.

3.7. Распределение ресурсов

Проблему оптимального таргетирования можно рассматривать как проблему распределения ресурсов, когда некоторые ограниченные ресурсы, такие как купоны, должны распределяться между клиентами. До сих пор мы концентрировались только на этом типе распределения и игнорировали тот факт, что маркетинговые действия необходимы для принятия решений о распределении многих других ресурсов. Корпоративная стратегия маркетинга обычно включает принятие решений о распределении ресурсов между маркетинговой и немаркетинговой деятельностью, товарами, стадиями жизненного цикла продукции, рынками и территориями, бизнес-целями, маркетинговыми каналами и типами коммуникаций [Carpenter and Shankar, 2013]. Многие из этих решений, такие как распределение ресурсов между

маркетинговой и исследовательской деятельностью, носят стратегический характер и поэтому не могут рассматриваться программной системой. Другие решения имеют тактический характер, и система может включать модули, автоматизирующие или, по крайней мере, облегчающие процесс принятия решений. Вместе с тем следует иметь в виду, что таргетирование является одной из наиболее тактических и технических задач распределения ресурсов, и автоматизация других решений о распределении ресурсов становится еще более сложной задачей.

Моделирование и оптимизация распределения ресурсов между маркетинговыми мероприятиями и возможностями называется *моделированием маркетингового комплекса* (Marketing Mix Modeling, MMM). Его можно рассматривать как статистический анализ влияния разных компонентов комплекса, таких как рекламные акции и цены, на показатели эффективности бизнеса — продажи и доходы. В этом разделе мы сосредоточимся на двух задачах распределения ресурсов: распределение между каналами и распределение между бизнес-целями, а также обсудим, как эти задачи решаются с помощью методов MMM.

3.7.1. Распределение между каналами

Система таргетирования часто имеет в своем распоряжении несколько маркетинговых каналов, и каждый канал имеет свою структуру затрат, аудиторию и эффективность. Канал прямой почтовой рассылки, например, может иметь гораздо более высокую стоимость сообщения, чем канал электронной почты, но способен обеспечить более высокий процент отклика для определенных категорий клиентов. Это требует оптимизации маркетинговых коммуникаций в отношении выбора канала. Один из возможных подходов к решению этой задачи — оптимизация на уровне клиента, когда канал выбирается с помощью модели отклика, учитывающей вероятности и затраты для конкретного канала. Другой подход заключается в оптимизации распределения глобального бюджета между каналами для максимизации доходов. Это иногда называют *моделированием комплекса каналов* (channel mix modeling). Эти две методологии можно рассматривать, соответственно, как восходящие и нисходящие решения, и обе они имеют большое значение.

Моделирование комплекса каналов — это набор методов статистического анализа, фокусирующихся на следующих описательных и прогнозирующих вопросах:

- Какой процент дохода (или другой меры качества работы) зависит от каждого канала или типа коммуникации?
- Как увеличение или уменьшение расходов для данного канала повлияет на доход?
- Каково оптимальное распределение бюджета между каналами?

Интуитивно можно ожидать, что ответы на эти вопросы сможет дать регрессионная модель, выражающая интересующий показатель как функцию от активности канала. Проблема, однако, в том, что зависимость между активностью и наблюдаемым показателем может быть сложной по нескольким причинам. Во-первых, измерить активность канала напрямую можно только как текущее количество электронных писем или показов онлайн-рекламы, но отклики клиентов обычно задерживаются и растянуты во времени. Во-вторых, одновременно может проводиться сразу несколько кампаний, а мы можем наблюдать только кумулятивный эффект. Наконец, зависимость между активностью канала и величиной отклика часто нелинейна из-за эффекта насыщения. Одной из популярных моделей комплекса каналов, учитывающей эти эффекты, является *модель Adstock*¹ [Broadbent, 1979].

Модель Adstock основана на предположении, что каждый данный период продаж отчасти подвержен сохранившемуся влиянию предыдущей рекламы. Допустив, что на данный момент у нас есть только один рекламный канал, обозначим активность канала, измеренную в потраченных рублях или в количестве сообщений за период времени t как x_t , интересующую бизнес-показатель, обычно объем продаж или доход, как y_t , и текущий эффект, вызванный воздействием на бизнес-показатель, как a_t . Переменная эффекта a_t называется Adstock. Согласно предположениям, модель Adstock можно выразить следующим образом:

$$a_t = x_t + \lambda \cdot a_{t-1}, \quad (3.41)$$

где λ — параметр затухания, соответствующий доле эффекта, сохранившейся к данному периоду времени. Например, величина 0,4 параметра означает, что воздействие в одном периоде сохранит 40 % влияния в следующем периоде. Иначе говоря, модель Adstock предполагает, что каждое новое маркетинговое действие поднимает лояльность и ориентированность на новый уровень, но затем лояльность постепенно затухает, пока вновь не будет поднята следующей порцией действий. Разложив рекурсивное уравнение 3.41, получим:

$$a_t = x_t + \lambda x_{t-1} + \lambda^2 x_{t-2} + \lambda^3 x_{t-3} + \dots \quad (3.42)$$

Обратите внимание, что, по сути, это сглаживающий фильтр, применяемый к входной последовательности. На практике можно смело предположить, что эффект воздействия конечен и ограничен n периодами, поэтому перепишем преобразование Adstock исходной последовательности как

$$a_t = x_t + \sum_{j=1}^n \lambda^j \cdot x_{t-j}. \quad (3.43)$$

¹ Модель забывания. — Примеч. пер.

В этом случае наблюдаемый бизнес-показатель можно оценить как линейную функцию от сохранившегося влияния (adstock):

$$\hat{y}_t = w a_t + c, \quad (3.44)$$

где w — весовой коэффициент, а c — базовое значение при отсутствии сохранившегося влияния (adstock). В случае нескольких каналов предполагается, что сохранившееся влияние (adstock) имеет аддитивный характер, поэтому полное определение модели является линейной регрессией по сохранившимся влияниям:

$$\hat{y}_t = \sum_{i=1}^n w_i a_{it} + c, \quad (3.45)$$

при этом каждый канал моделируется со своим параметром затухания λ_i , из-за чего полная модель требует оценки базового параметра c , n параметров затухания λ_i и весовых коэффициентов w_i для n каналов. Мы можем обучить модель, решив следующую задачу для наблюдаемых выборок y_t :

$$\min_{c, w, \lambda} \sum_t |y_t - \hat{y}_t|^2. \quad (3.46)$$

Обученная модель позволяет оценить влияние увеличения или уменьшения бюджетов для каналов и относительный вклад каждого канала в целевой показатель:

$$z_{it} = \frac{w_i a_{it}}{\sum_j w_j a_{jt}}. \quad (3.47)$$

Это значение можно усреднить по времени и получить средний относительный вклад канала. Эффективность канала можно оценить как отношение между абсолютным вкладом и бюджетом канала, то есть количеством проданных единиц на каждый рубль, потраченный на маркетинговое воздействие через этот канал. Следующий пример иллюстрирует, как создать и использовать модель Adstock.

ПРИМЕР 3.6

Рассмотрим ретейлера, использующего для рассылки рекламы два канала: электронную почту и СМС. Ретейлер может измерять и контролировать интенсивность маркетинговых коммуникаций по каждому из каналов, устанавливая бюджет и ограничивающие правила. Ретейлер также следит за объемом продаж. На рис. 3.23 построен график изменения этих метрик в течение 20 последовательных временных интервалов (опустим таблицу с числовыми значениями для экономии места).

Модель Adstock можно обучить, решив уравнение 3.46 численными методами оптимизации. Установив протяженность n окна затухания равной 3, получим следующие оценки для параметров модели:

базовый уровень: $c = 28,028$

e-mail: $\lambda_{email} = 0,790$ $w_{email} = 1,863$

СМС: $\lambda_{sms} = 0,482$ $w_{sms} = 4,884$.

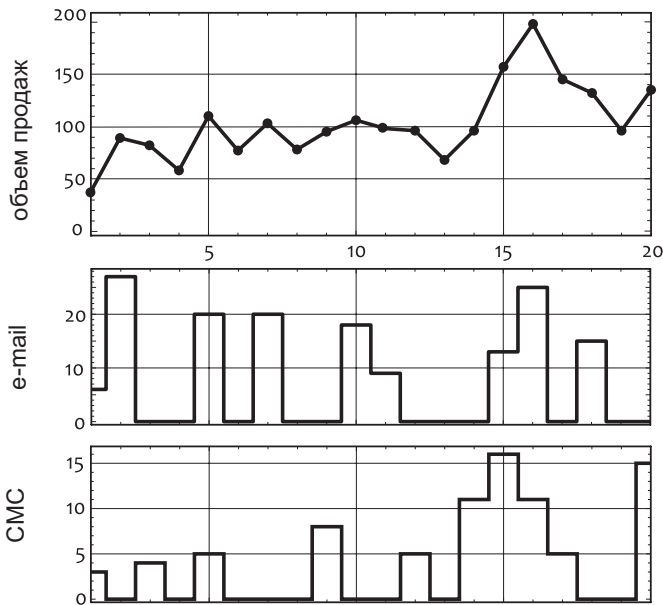


Рис. 3.23. Данные для моделирования сохранившегося влияния (adstock): объем продаж, интенсивность рассылки по электронной почте и СМС

Эти параметры можно использовать для расчета сохранившегося влияния (adstock), как определено в выражении 3.43. Структуру объема продаж можно визуализировать в виде трех слоев, уложенных друг на друга: базовый объем продаж определяется константой c , вклад, обусловленный электронной почтой, оценивается как сохранившееся влияние электронной почты, масштабированное коэффициентом w_{email} , а вклад СМС — как сохранившееся влияние СМС, взвешенное коэффициентом w_{sms} . Эта структура соответствует выражению 3.45, а результат показан на рис. 3.24. Такая декомпозиция позволяет оценить эффективность каждого канала и оптимизировать распределение бюджета.

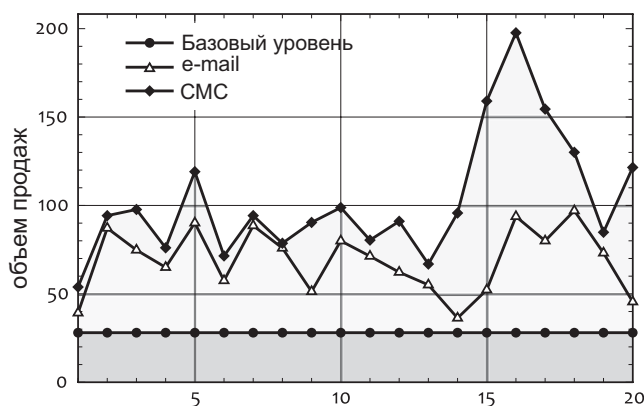


Рис. 3.24. Разложение объема продаж на слои, описывающие вклад разных маркетинговых каналов

Основная модель Adstock учитывает перекрывающиеся маркетинговые воздействия и эффект затухания, но не учитывает рекламную насыщенность, о которой говорилось выше. В целом увеличение интенсивности воздействия увеличивает охват кампании и, соответственно, увеличивает спрос. Однако зависимость между интенсивностью и спросом нелинейна. Как правило, она следует *закону убывания доходности*, поэтому в какой-то момент увеличение затрат на маркетинговую деятельность дает более низкий прирост спроса. Модель Adstock может учесть этот эффект насыщения нелинейным преобразованием переменной интенсивности. Для этой цели, например, можно использовать сигмоидную (логистическую) функцию, в результате чего рекурсивное уравнение 3.41 можно переписать следующим образом:

$$a_t = \frac{1}{1 + \exp(-\mu \cdot x_t)} + \lambda \cdot a_{t-1}, \quad (3.48)$$

где μ — дополнительный параметр модели, управляющий крутизной логистической кривой. Модель Adstock можно дополнить или изменить разными способами, чтобы учесть дополнительные эффекты, встречающиеся на практике. Например, часто бывает нужно учесть сезонность спроса, для чего в модель можно добавить дополнительные переменные. В данный момент, моделируя комплекс каналов, можно использовать преимущества методов моделирования спроса, которые подробно рассматриваются в главе 6. Другой пример — предположение о геометрическом запаздывании, используемое в модели Adstock, является несколько

ограничивающим, потому что запаздывание по времени может иметь более сложную форму. Фактически модель, описываемая уравнениями 3.42 и 3.43, известна в эконометрике как модель распределенного запаздывания Койка (Koyck), которая относится к семейству моделей распределенного запаздывания [Koyck, 1954]. Это семейство включает ряд более гибких альтернатив, в том числе полиномиальную модель распределенного запаздывания, более гибкую и простую для аппроксимации, чем модель Койка [Almon, 1965; Hall, 1967].

3.7.2. Распределение по целям

Для оптимизации таргетирования программная система может использовать такие цели, как рост LTV или привлечение-максимизацию-удержание. Рентабельность вложений во всех этих случаях можно оценить как увеличение LTV или непосредственное увеличение чистой прибыли, в соответствии с реализацией моделирования отклика. Выбор между этими целями и глобальной оптимизацией окупаемости вложений является стратегическим вопросом, который необязательно должен решаться программной системой. Тем не менее система может дать некоторые рекомендации, как лучше распределить бюджет между целями, чтобы максимизировать общую рентабельность [Blattberg and Deighton, 1996].

Из раздела 3.5.7 мы узнали, что основным фактором, влияющим на LTV, является коэффициент удержания, поэтому LTV можно рассматривать как функцию от коэффициента удержания. Можно предположить, что коэффициент удержания, в свою очередь, является функцией маркетингового бюджета, потраченного на удержание. Например, зависимость между бюджетом и коэффициентом можно смоделировать следующим образом:

$$r = r_{\max} \left(1 - e^{-k_r R}\right), \quad (3.49)$$

где R — бюджет на удержание одного клиента, r_{\max} — оценка максимальной степени удержания (потолок), которой можно достичь при неограниченном бюджете, а k_r — коэффициент, определяющий скорость приближения доли к потолку. Аналогично коэффициент привлечения a (доля откликов на кампанию привлечения) можно смоделировать как функцию от бюджета привлечения:

$$a = a_{\max} \left(1 - e^{-k_a A}\right), \quad (3.50)$$

где A — бюджет привлечения одного клиента, a_{\max} — оценка максимальной доли отклика, а k_a — параметр, определяющий чувствительность изменения доли к изменению бюджета. Следовательно, чистую прибыль от привлечения данного клиента можно определить как

$$a \cdot \text{LTV}(r) - c, \quad (3.51)$$

где c — стоимость привлечения перспективного клиента. Тогда общую задачу оптимизации бюджетов A и R можно определить следующим образом:

$$\begin{aligned} \max_{A, R} \quad & N_p (a \cdot \text{LTV}(r) - c) + N_c \cdot \text{LTV}(r), \\ \text{с условием} \quad & A + R \leq \text{общий бюджет}, \end{aligned} \quad (3.52)$$

где N_p — общее количество потенциальных клиентов, а N_c — общее количество текущих клиентов. Первый член уравнения 3.52 соответствует доходам от новых клиентов, а второй — доходам от существующих клиентов, поэтому фактически это задача оптимизации доходов. Уравнение 3.52 определяет задачу оптимизации в терминах агрегированных и усредненных значений, но его можно легко переписать как сумму отдельных LTV по всем клиентам, чтобы получить более точные оценки с помощью предиктивных моделей.

3.8. Онлайн-реклама

Принципы таргетирования продвижения, рассмотренные в предыдущих разделах, ориентированы на товары широкого спроса и традиционную розничную среду. Очевидно, что многие из этих принципов применимы к другим маркетинговым средам, однако их реализация в значительной степени зависит от имеющихся данных и точного определения бизнес-целей, которые могут различаться в разных средах. Продолжим анализ онлайн-рекламы, которая, пожалуй, является наиболее важным и хорошо развитым приложением алгоритмического маркетинга и является отличным примером среды, где техническая инфраструктура и потоки данных настолько сложны, что бизнес-цели невозможно понять и достигнуть без тщательного изучения технических возможностей и ограничений.

3.8.1. Среда

Среда онлайн-рекламы очень сложная и многообразная, потому что представляет рынок, где тысячи компаний продают и покупают рекламные ресурсы, предлагают и используют технические системы, автоматизирующие процесс покупки, а также контролируют и измеряют качество и эффективность рекламных кампаний. Дополнительная сложность проистекает из того факта, что хотя большая часть терминологии и стандартных предложений, как правило, остается общей для всей отрасли, существует множество вариаций и сквозных решений, которые появляются по мере развития отрасли. Высокая сложность экосистемы онлайн-рекламы

затрудняет учет всех важных аспектов среды в одном представлении, поэтому для поддержки обсуждения экономических и бизнес-целей начнем с упрощенной модели, изображенной на рис. 3.25.

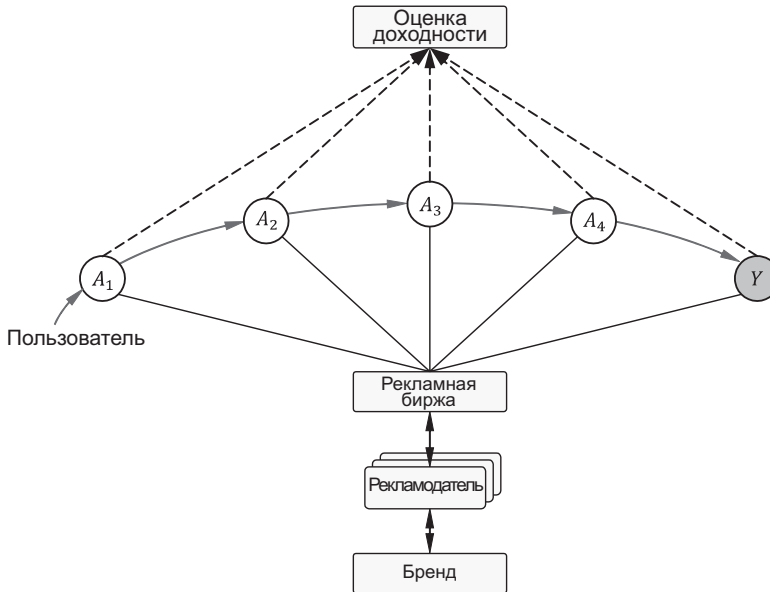


Рис. 3.25. Среда онлайн-рекламы

На рис. 3.25 показаны отношения между следующими ключевыми сущностями, составляющими ландшафт онлайн-рекламы:

- *Бренд*, также известный как *маркетолог*, является продавцом товаров или услуг. Бренд инвестирует деньги в рекламные кампании и рассчитывает получить отдачу в виде увеличения продаж и улучшения взаимоотношений с клиентами.
- *Рекламодатель*, или агентство, которое проводит рекламные кампании от имени бренда. Рекламодатель, как правило, пытается достичь тех же целей, что и бренд, но выбираемая им стратегия зависит от модели оплаты, установленной между брендом и рекламодателем, а также методологии измерения эффективности кампании. Бренд может работать с несколькими агентствами, конкурирующими друг с другом в рамках одной кампании.
- Рекламодатели могут выходить на интернет-пользователя, являющегося текущим или потенциальным клиентом бренда, по разным *каналам*. Примерами каналов могут служить рекламные баннеры на веб-странице, оплаченные

результаты на странице, возвращаемой поисковиком, и онлайн-видеореклама. В общем случае набор каналов не ограничивается интернет-каналами и может включать другие средства массовой информации, такие как телевизионная реклама или печатные каталоги.

- Каждый канал представлен несколькими *издателями*, например веб-сайтами. Издатели продают рекламные площади, то есть доступные слоты, которые могут содержать рекламные объявления.
- Издатели и рекламодатели связаны через рекламную биржу. При наличии свободной рекламной площади издатель посылает бирже запрос на размещение объявлений (например, когда пользователь открывает веб-страницу), а биржа распределяет запросы между рекламодателями, которые, в свою очередь, могут купить доступный рекламный слот и показать объявление пользователю. Биржа часто организована как аукцион Викри (аукцион второй цены), который обрабатывает каждый запрос объявления в режиме реального времени, поэтому биржу обычно называют процессом *торгов в режиме реального времени* (Real-Time Bidding, RTB).
- *Пользователь* является получателем рекламы, доставляемой по каналам. Пользователь может взаимодействовать с несколькими каналами и издателями, получая реализации объявлений, известные как *показы*. С точки зрения бренда пользователь в конечном итоге *конвертируется* в направлении некоторого желаемого результата, такого, как покупка на веб-сайте бренда, или не конвертируется. Следовательно, для каждого пользователя существует воронка последовательных показов событий A , которая заканчивается результатом Y , как показано на рис. 3.25.
- Наконец, показы и конверсия отслеживаются *системой оценки доходности*. Мы рассматриваем эту систему как абстрактную сущность, которая может отслеживать личность пользователя по каналам и издателям и вести учет, какое впечатление получил пользователь от конкретного рекламодателя в конкретный момент времени. Целью системы оценки доходности является измерение эффективности рекламной кампании и анализ вклада отдельных каналов, рекламодателей и сегментов пользователей. Система оценки обычно собирает информацию с помощью *пикселей* отслеживания, прикрепляемых к рекламным баннерам и веб-страницам; пользователи идентифицируются посредством файлов cookie веб-браузера. Однако процесс оценки может использовать дополнительные источники данных, такие как покупки в обычных магазинах, сопоставлять эти данные с онлайн-профилями посредством идентификаторов карт постоянного клиента или кредитных карт и измерять причинно-следственные эффекты между онлайн- и офлайн-каналами.

В модели среды, приведенной выше, бренд полагается на систему оценки доходности для измерения эффективности отдельных рекламодателей и рекламных кампаний в целом. Метрики, полученные системой оценки, напрямую преобразуются в гонорары рекламодателей, затраты и доходы бренда, поэтому в следующем разделе мы рассмотрим модели оценки и их влияние на выбор стратегии рекламодателем.

3.8.2. Цели и оценка

По аналогии с промоакциями, бизнес-целями бренда движет желание перевести отношения с потребителями на более высокий уровень:

УЗНАВАЕМОСТЬ БРЕНДА. Маркетолог, как правило, заинтересован в том, чтобы сделать свой бренд узнаваемым для потенциальных клиентов и связанным с определенной категорией продуктов, такой как безалкогольные напитки или автомобили класса люкс, даже если это не сразу приводит к их конверсии.

ПРИВЛЕЧЕНИЕ КЛИЕНТОВ. Цель привлечения — заинтересовать потенциальных клиентов, не взаимодействующих с брендом, и привести их к конверсии.

РЕТАРГЕТИРОВАНИЕ. Ретаргетирование, также известное как ремаркетинг, направлено на перспективных клиентов, которые уже взаимодействовали с брендом, когда есть потенциал для развития отношений с ними. Типичным примером могут служить интернет-пользователи, которые посещали сайт бренда один или несколько раз, но не были конвертированы.

К этим основным целям можно добавить дополнительные ограничения, важные для бренда. Например, бренд может не захотеть размещать рекламу на веб-сайтах со взрослым или агрессивным контентом или с контентом, воспитывающим ненависть. В идеале договор между брендом и рекламодателем должен быть оформлен так, чтобы рекламодателю платили за достижение вышеуказанных целей. Более конкретно желаемые свойства контракта можно описать следующим образом:

- Процессы таргетирования и проведения торгов должны определяться бизнес-целями кампании (например, узнаваемость бренда или ретаргетирование) и ограничиваться дополнительными правилами, такими как *репутационная безопасность бренда*.
- Эффект кампании должен быть измеримым, а показатели должны точно отражать добавленную рекламодателем стоимость. Иначе говоря, метрики должны отвечать на вопрос: что случится с бизнес-целью, если убрать рекламодателя. Обратите внимание, что этот вопрос напрямую связан с моделированием увеличения эффективности, которое рассматривалось выше в этой главе.

- Должна иметься возможность получить ответ на вопрос об удалении рекламодателя в случае, если несколько рекламодателей работают на один и тот же бренд. Суммы оплаты должны начисляться рекламодателям пропорционально их вкладу в общий прирост стоимости.

К сожалению, сложно определить контракт, полностью отвечающий вышеуказанным критериям. Бизнес-цели могут быть формализованы по-разному, и оценка добавочной стоимости также является нетривиальной статистической и технической задачей. Давайте сначала познакомимся с некоторыми основными методами, широко используемыми на практике, а затем в последующих разделах обсудим нерешенные вопросы и ограничения.

С точки зрения бренда, общую эффективность кампании можно оценить по *стоимости привлечения* (Cost Per Acquisition, CPA), которая определяется как общая стоимость C_{camp} кампании, деленная на общее количество конверсий N_{conv} :

$$CPA = \frac{C_{camp}}{N_{conv}}. \quad (3.53)$$

Конверсию можно определить по-разному. Один из возможных подходов — подсчет *действий после просмотра*, то есть подсчет пользователей, посетивших сайт бренда или совершивших покупку в течение определенного периода времени (например, в течение недели) после показа рекламы. Более простой метод — подсчет немедленных щелчков на объявлениях, который называется моделью *затрат на клик* (Cost Per Click, CPC). С точки зрения рекламодателя имеет смысл представить стоимость кампании как произведение количества показов на среднюю стоимость одного показа, то есть показатель CPA можно выразить следующим образом:

$$CPA = \frac{N_{impr} \cdot \mathbb{E}[c_{impr}]}{N_{conv}} = \frac{1}{CR} \cdot \mathbb{E}[c_{impr}], \quad (3.54)$$

где N_{impr} — общее количество показов в рамках кампании, $\mathbb{E}[c_{impr}]$ — средняя цена, уплаченная брендом за один показ, а CR — коэффициент конверсии. Доход рекламодателя — разница между ценой, уплаченной брендом, и стоимостью ставки на торгах в реальном времени, поэтому эквивалент стоимости привлечения для рекламодателя можно определить следующим образом:

$$CPA_a = \frac{1}{CR} \cdot \mathbb{E}[c_{impr} - c_{bid}]. \quad (3.55)$$

Нам нужно определить контракты для c_{impr} и c_{bid} , чтобы оценить приведенные выше выражения для CRA и CRA_a . Обычно бренд платит фиксированную цену c_{impr} , хотя на практике используются два разных типа контрактов:

- *Стоимость действия* (Cost Per Action), также известная как *стоимость привлечения* (Cost Per Acquisition, CPA) или *плата за привлечение* (Pay Per Acquisition, PPA). Бренд платит фиксированную комиссию за каждую конверсию, измеряемую системой оценки.
- *Стоимость тысячи показов* (Cost Per Mile, CPM). Бренд платит фиксированную плату за каждый показ, но в конечном итоге измеряет общую стоимость привлечения с помощью системы оценки.

Оба подхода эквивалентны в том смысле, что рекламодатель должен минимизировать показатель CPA, чтобы удовлетворить клиента, даже для CPM-контрактов. Фиксированная плата подразумевает возможность оптимизации метрики CPA в уравнении 3.54 путем максимизации коэффициента конверсии CR . Однако значение c_{bid} в уравнении 3.55 не является фиксированным и напрямую влияет на коэффициент конверсии, поэтому оптимизация метрики CRA_a требует совместной оптимизации CR и c_{bid} .

Последняя область, которую мы должны охватить, — оценка доходности в случае нескольких рекламодателей. Самый простой подход — *оценка последнего касания* (Last-Touch attribution, LT), согласно которому вся оплата отдается последнему показу, предшествовавшему конверсии. Соответственно, цель рекламодателя в рамках модели LT — идентифицировать клиентов, которые, скорее всего, конвертируются сразу после показа.

Параметры CPA и LT — будем называть их моделью CPA-LT — дают достаточно полное и формальное определение задачи, которое можно использовать для оптимизации процесса таргетирования. Однако модель CPA-LT является чрезмерно упрощенной и имеет ряд проблем и ограничений:

- Отсутствует явная связь с бизнес-целью. Модель не различает цели, описанные выше: привлечение, узнаваемость и ретаргетирование. В действительности принципы CPA-LT ориентированы на потребителей с высокой предрасположенностью к покупке, что подразумевает большой уклон в сторону ретаргетирования и тактического привлечения, а не узнаваемости и стратегического привлечения.
- Модель предполагает оптимизацию отклика, а не подъем. При определенных обстоятельствах это может привести к бессмысленным результатам. Например, метод таргетирования, выявляющий только пользователей с высокой предрасположенностью к конверсии без показов, будет иметь очень высокую

эффективность в рамках модели CPA-LT, хотя это вряд ли будет хорошим подходом с точки зрения рентабельности инвестиций.

- Оценка доходности методом «последнего касания» побуждает рекламодателей к обману и паразитированию на усилиях друг друга. Например, рекламодатель может купить много низкокачественных рекламных площадей, допустим, внизу веб-страниц, чтобы «затронуть» как можно больше пользователей (так называемый метод *ковровых бомбардировок*).

В следующем разделе мы обсудим способы оптимизации стратегии таргетирования и торгов в рамках модели CPA-LT, а затем рассмотрим способы устранения недостатков этой модели с помощью более сложных оценок доходности и контролируемых экспериментов.

3.8.3. Таргетирование для модели CPA-LT

Основная цель таргетирования в модели CPA-LT — выявление пользователей, которые с высокой вероятностью конвертируются вскоре после показа. Как и в случае с таргетированием продвижения, для решения этой задачи используем вариант моделирования методом подобию, при этом мы должны явно учесть информацию о реакции пользователя на рекламу, а не выбирать естественных покупателей, опираясь на историю покупок. В частности, мы должны учесть эффективность текущей рекламы, а это, в свою очередь, означает необходимость динамически корректировать метод таргетирования, исходя из наблюдаемых результатов. Другими словами, мы должны создать самонастраивающийся метод таргетирования.

Предположим, что рекламодатель может получить следующие данные из профилей потребителей:

- Посещавшиеся URL-адреса. Рекламодатель анализирует запросы на предложения и другие источники данных партнеров, позволяющие выяснить историю просмотров пользователя. URL-адреса могут включать домены, например google.com, и адреса конкретных страниц.
- Атрибуты пользователя. Вместе с URL-адресами рекламодатель может получить дополнительную информацию о пользователе, например, характеристики устройства и приложений, с помощью которых пользователь просматривает страницы, географическое положение, время, проведенное на странице, и некоторые другие сведения.
- Ставки и показы. Рекламодатель может отслеживать ставки, которые он сделал для данного пользователя, и показы рекламы ему.

- Щелчки. Рекламодатель может получить информацию о взаимодействии пользователя с показанной ему рекламой.
- Конверсии. Бренд может предоставить рекламодателю информацию о конверсиях на своем веб-сайте.
- Дополнительные данные бренда. Бренд может предоставить дополнительные данные, например, какие товары пользователь просматривал на веб-сайте бренда.

На основе этих данных можно разработать признаки для предиктивного моделирования. Известно, что посещавшиеся URL-адреса и производные характеристики, такие как периодичность и частота, содержат много предиктивной информации о конверсии. Однако основная проблема заключается в большом количестве наблюдаемых URL-адресов, из-за чего модель, использующая вектор с двоичными элементами, указывающими, посещал ли пользователь тот или иной URL-адрес, может иметь миллионы измерений.

Самый простой подход к задаче самонастраивающегося таргетирования заключается в том, чтобы запустить кампанию со случайным таргетированием, то есть сделать ставку на случайных людей, затем дождаться достаточного количества конверсий и обучить модель количественной оценки, используя конвертированных пользователей в качестве положительных примеров, и не конвертированных — в качестве отрицательных, как показано на рис. 3.26. Однако это не самый оптимальный подход, потому что события конверсии очень редки в случайном



Рис. 3.26. Желательная выборка для задачи таргетирования. Заштрихованные круги соответствуют положительным и отрицательным примерам

таргетировании, а размерность профилей пользователей, как говорилось выше, очень высока, поэтому создание достаточного набора обучающих данных с использованием случайных ставок в начале кампании может оказаться неприемлемо и непрактично дорогостоящим [Dalessandro et al., 2012a].

Существует много разных методов улучшения базового подхода, описанного выше. В остальной части этого раздела мы познакомимся с поэтапной методологией таргетирования, описанной в работах Dalessandro et al., 2012a и Perlich et al., 2013, которая предлагает исчерпывающее практическое решение для самонастраивающегося таргетирования. Суть ее состоит в том, чтобы выполнить процесс таргетирования в три последовательных этапа: рассчитать близость бренда, включить отклик на объявление и оценить качество рекламной площади, рассчитав сумму ставки.

3.8.3.1. Близость бренда

Цель этого этапа — оценить вероятность конверсии Y независимо от влияния рекламы, то есть рассчитать безусловную *близость бренда* $\Pr(Y|u)$ для пользователя u . Если историческая информация о посетителях сайта бренда доступна до начала кампании, рекламодатель может создать модели пользователей, которых можно конвертировать, взяв профили конвертированных пользователей в качестве положительных примеров и случайные профили пользователей интернета в качестве отрицательных. Обратите внимание, что эта выборка отличается от желаемой, изображенной на рис. 3.26. Этот шаг, по сути, похож на применение метода подобию с использованием посещавшихся URL-адресов в качестве признаков и конверсий в качестве меток для моделирования безусловной близости бренда:

$$\begin{aligned}\phi(u) &= \Pr(Y|u) = \\ &= \Pr(Y|URL_1, \dots, URL_n),\end{aligned}\tag{3.56}$$

где URL_i — бинарные метки, равные единице, если пользователь посетил соответствующий URL, и ноль в противном случае. Рекламодатель может использовать различные определения URL и конверсии для построения нескольких моделей $\phi_{u1}, \dots, \phi_{uk}$, использующих разные индикаторы близости:

- URL-адреса могут объединяться в кластеры, а метки URL_i — заменяться бинарными метками кластеров, которые указывают, посетил ли пользователь некоторый URL-адрес из данного кластера. Это уменьшает размерность задачи, что может быть полезно, если число доступных событий конверсии относительно невелико. Расстояние между URL-адресами, необходимое для класте-

ризации, можно определять на основе оценок качества рекламных площадей, которые рассматриваются далее в этом разделе.

- Конверсию можно определить как посещение сайта бренда, покупку после показа рекламы или любую покупку.

Модель близости бренда можно использовать для оценки пользователей в начале кампании, когда фактические данные об откликах на рекламные объявления еще не доступны. Следующий шаг — добавление новых данных, когда они станут доступны, и корректировка результатов.

3.8.3.2. Моделирование откликов на рекламные объявления

Целью этапа моделирования отклика является оценка условной вероятности конверсии $\Pr(Y|u, a)$ для рекламного объявления a . В основном этот шаг делает все то же самое, что и базовый подход, описанный в начале этого раздела, — рекламодатель использует модель близости ϕ для таргетирования целевой аудитории в начале кампании, но при этом объявления показываются также небольшому количеству случайных людей, чтобы получить желательную выборку, как показано на рис. 3.26. Но, в отличие от базового метода, теперь вместо исходных URL-адресов в качестве признаков можно использовать выходные данные предыдущего этапа, что делает процесс обучения более эффективным. Формулу близости бренда можно дополнить признаками с информацией о пользователях f_{u1}, \dots, f_{ur} , такими как тип браузера и географическое местоположение, поэтому модель можно описать следующим образом:

$$\begin{aligned}\psi_a(u) &= \Pr(Y|u, a) = \\ &= \Pr(Y|\phi_{u1}, \dots, \phi_{uk}, f_{u1}, \dots, f_{ur}).\end{aligned}\tag{3.57}$$

Ключевое различие между моделями безусловной близости ϕ и предрасположенности к конверсии ψ заключается в выборке: семейство моделей ϕ классифицирует пользователей на конвертировавшихся или неконвертировавшихся, независимо от рекламы, тогда как модель ψ классифицирует пользователей на откликнувшихся и неоткликнувшихся и зависит от рекламы. Однако оценки, полученные с помощью моделей ϕ , обладают высокой способностью предсказывать отклик, обеспечивая разумные начальные значения для ψ и делая переоценку ψ более эффективной по мере поступления фактических данных об откликах.

3.8.3.3. Качество рекламной площади и ставок

Заключительный этап — включение дополнительной информации, отсутствующей в оценках, полученных моделью ψ , и определение фактической цены ставки для предложения на рекламной бирже. В предположении, что рекламная биржа

реализует аукцион второй цены, оптимальную цену ставки можно рассчитать как ожидаемую ценность конверсии $v(Y)$:

$$b_{opt} = \mathbb{E}[v(Y)] = \Pr(Y | u, a) \cdot v(Y). \quad (3.58)$$

Ценность конверсии $v(Y)$ обычно можно считать постоянной для всех пользователей и включающей некоторую базовую цену ставки b_{base} , установленную рекламодателем и зависящую от контракта с брендом и особенностей биржи. Следовательно, оценки предрасположенности можно рассматривать как множители, масштабирующие базовую цену.

Оценки предрасположенности, возвращаемой моделью ψ , обычно достаточно для таргетирования и принятия решения относительно размера ставки. Цену ставки для данного пользователя можно вычислить как

$$b(u) = b_{base} \cdot s_1(\psi_a(u)), \quad (3.59)$$

где $s_1(\cdot)$ — некоторая масштабирующая функция для оценки ψ . В частности, $s_1(\cdot)$ может отображать все оценки ниже некоторого порога в ноль (нет ставки), а порог определить, исходя из желаемого количества показов и других соображений, как рассматривалось выше в контексте рекламных акций.

Процесс таргетирования, как описывалось до сих пор, учитывает профили пользователей и рекламные объявления, но не контекст показов, то есть рекламную площадь. Качество рекламной площади играет важную роль по нескольким причинам [Perlich et al., 2013]:

- Рекламная площадь несет информацию о намерении пользователя совершить покупку и релевантности объявления для пользователя. Например, реклама отеля будет иметь более высокий коэффициент конверсии на туристических сайтах, чем на новостных.
- Восприятие рекламы зависит от контекста. Например, пользователи, читающие сложные технические материалы, могут уделять рекламе меньше внимания, чем посетители развлекательных сайтов; некоторые рекламные слоты могут быть плохо расположены, и пользователь должен прокрутить страницу вниз, чтобы увидеть их, и т. д.

Следовательно, рекламодатель может рассчитывать на лучшие результаты, используя вероятность $w_a(u, i) = \Pr(Y | u, a, i)$, где i — рекламная площадь. Соотношение $w_a(u, i)$ и его ожидание по всем площадям $w_a(u) = \mathbb{E}_i[w_a(u, i)]$ можно использовать в роли меры качества площади, потому что оно показывает, насколько лучше или хуже площадь i по сравнению со средней площадью. Эту метрику можно использовать как дополнительный множитель для масштабирования ставки:

$$b(u) = b_{base} \cdot s_1(\psi_a(u)) \cdot s_2\left(\frac{w_a(u, i)}{w_a(u)}\right). \quad (3.60)$$

Обратите внимание: даже при том, что использованные нами обозначения подразумевают равенство $w_a(u)$ и $\psi_a(u)$, рекламодатель может использовать разные выборки данных и модели для оценки w и ψ , в зависимости от имеющихся данных и других факторов. Крутизна функций масштабирования $s_1(\cdot)$ и $s_2(\cdot)$ определяет компромисс между коэффициентами конверсии и CPA рекламодателя. Крутые функции масштабирования (например, возвращающие ноль, если аргумент ниже порога, и очень высокое значение в противном случае) обычно максимизируют коэффициент конверсии, но могут быть неоптимальными с точки зрения CPA. Функции масштабирования, близкие к функции тождества, оптимизируют CPA, как следует из теоретического уравнения 3.58, но могут быть неоптимальными с точки зрения коэффициентов конверсии.

3.8.4. Оценка для случая с несколькими каналами

Очевидным ограничением оценки последнего касания является игнорирование усилий, предшествовавших последнему показу. Обойти его можно, используя более сложные методы оценки, которые распределяют затраты в соответствии с положением рекламодателя в воронке. Несколько примеров таких статических моделей приводится в табл. 3.9. Однако статическая оценка по весу не позволяет оценить вклад отдельных рекламодателей в общий эффект кампании. Нам нужно создать *алгоритмический метод оценки*, который измерит фактический вклад и позволит бренду вознаградить лучших рекламодателей или каналы и избавиться от худших из них.

Таблица 3.9. Модель статической оценки. В таблице приводятся процентные доли, присвоенные каждому из пяти показов A_1, \dots, A_5

Модель	A_1	A_2	A_3	A_4	A_5
Первый показ	100 %	—	—	—	—
Первый щелчок	—	100 %	—	—	—
Последнее касание	—	—	—	—	100 %
Линейная	20 %	20 %	20 %	20 %	20 %
Позиционная	35 %	10 %	10 %	10 %	35 %
Временной спад	10 %	15 %	20 %	25 %	30 %

Предположим, что бренд работает с сетью рекламодателей или каналами $C = \{C_1, \dots, C_n\}$. Эту сеть можно рассматривать как набор состояний, через которые может пройти пользователь перед конверсией, как показано на рис. 3.27. Причинно-следственный эффект канала C_k можно определить как разность между вероятностью конверсии для полного набора каналов и вероятностью конверсии при удалении канала C_k :

$$V_k = \Pr(Y | C) - \Pr(Y | C \setminus C_k). \quad (3.61)$$

Чтобы вычислить это выражение, можно перечислить все возможные подмножества из множества $C \setminus C_k$ и аппроксимировать причинно-следственный эффект для каждого подмножества в отдельности [Dalessandro et al., 2012b]:

$$V_k = \sum_{S \subseteq C \setminus C_k} w_{S,k} (\Pr(Y | S \cup C_k) - \Pr(Y | S)). \quad (3.62)$$

Коэффициенты $w_{S,k}$ моделируют распределение вероятностей конкретных реализаций S , то есть вероятность прохождения пользователем определенной последовательности каналов. Предположив равномерное распределение всех последовательностей, получаем:

$$w_{S,k} = \binom{|C| - 1}{|S|}^{-1} \cdot \frac{1}{|C|} = \frac{|S|! (|C| - 1 - |S|)!}{|C|!}, \quad (3.63)$$

потому что выстраиваем последовательности с длиной $|S|$ из множества $C \setminus C_k$ с кардинальностью $|C| - 1$. Например, причинно-следственный эффект канала C_3 в сети $C = \{C_1, C_2, C_3\}$ задается следующим уравнением:

$$\begin{aligned} V_3 = & \frac{1}{3} (\Pr(Y | C_1, C_2, C_3) - \Pr(Y | C_1, C_2)) + \\ & + \frac{1}{6} [(\Pr(Y | C_1, C_3) - \Pr(Y | C_1)) + (\Pr(Y | C_2, C_3) - \Pr(Y | C_2))] + \\ & + \frac{1}{3} (\Pr(Y | C_3) - \Pr(Y | \emptyset)). \end{aligned} \quad (3.64)$$

Формулу оценки 3.62 порой трудно реализовать на практике, поскольку длинные последовательности каналов имеют относительно низкую вероятность внедрения, что влияет на стабильность оценки [Dalessandro et al., 2012b; Shao and Li, 2011]. Возможно, разумнее отказаться от всех последовательностей S , включающих больше 2 каналов, чтобы получить более простую и стабильную модель [Shao и Li, 2011]:

$$\begin{aligned}
V_k^* &= \sum_{S \subseteq C \setminus C_k} w_{S,k} (\Pr(Y | S \cup C_k) - \Pr(Y | S)) = \\
&= w_0 [\Pr(Y | C_k) - \Pr(Y | \emptyset)] + \\
&\quad + w_1 \sum_{j \neq k} [\Pr(Y | C_j, C_k) - \Pr(Y | C_j)].
\end{aligned} \tag{3.65}$$

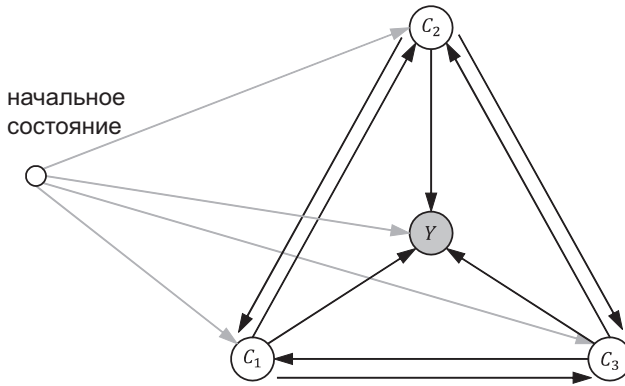


Рис. 3.27. Пример сети с тремя каналами

Также можно отбросить базовую вероятность конверсии $\Pr(Y | \emptyset)$, потому что она одинакова для всех каналов, а коэффициенты определены как

$$\begin{aligned}
w_0 &= \begin{pmatrix} |C| - 1 \\ 0 \end{pmatrix}^{-1} \frac{1}{|C|} = \frac{1}{|C|}; \\
w_1 &= \begin{pmatrix} |C| - 1 \\ 0 \end{pmatrix}^{-1} \frac{1}{|C|} = \frac{1}{(|C| - 1)|C|}.
\end{aligned} \tag{3.66}$$

Таким образом, причинно-следственный эффект можно выразить как

$$\begin{aligned}
V_k^* &= \frac{1}{|C|} \Pr(Y | C_k) + \\
&\quad + \frac{1}{(|C| - 1)|C|} \sum_{j \neq k} [\Pr(Y | C_j, C_k) - \Pr(Y | C_j)].
\end{aligned} \tag{3.67}$$

Вероятность конверсии $\Pr(Y | C_k)$ можно аппроксимировать как отношение числа конвертированных пользователей, прошедших через канал C_k , к общему числу пользователей, прошедших через канал. Вероятности второго порядка $\Pr(Y | C_j, C_k)$ для пары каналов можно аппроксимировать аналогично.

Уравнения 3.62 и 3.67 описывают практическое решение задачи получения оценки для случая с несколькими каналами. Однако стоит отметить, что существуют альтернативные решения. Например, можно построить регрессионную модель, предсказывающую конверсию на основе пройденных каналов, и сравнить величины коэффициентов регрессии [Shao and Li, 2011].

3.9. Оценка эффективности

Эффективность маркетинговых кампаний с трудом поддается оценке, потому что каждый потребитель имеет уникальные свойства, изменяющиеся с течением времени, и взаимодействует с брендом и маркетинговыми средствами информации по-своему, поэтому оценка любого положительного или отрицательного эффекта конкретного маркетингового действия всегда остается спорной. Маркетологи, как правило, не могут строго доказать эффективность действия, но могут попытаться организовать эксперимент или проанализировать собранные данные таким образом, чтобы рассматриваемые действия и их результаты были должным образом изолированы, и причинно-следственный эффект нельзя было отнести к внешним факторам. Это можно рассматривать как доказательство статистической значимости причинно-следственной связи между действиями и результатами.

Такая постановка задачи позволяет использовать развитый статистический аппарат, разработанный в других областях задолго до появления алгоритмического маркетинга. Важно отметить, что основа для экспериментов, созданная в таких областях, как биология и здравоохранение, специально адаптирована для работы со сценариями, структурно похожими на маркетинговые кампании.

3.9.1. Рандомизированные эксперименты

Рассмотрим простую маркетинговую кампанию, в ходе которой потенциальным клиентам рассылаются рекламные предложения с целью конвертировать их. Хотя наша конечная цель состоит в том, чтобы оценить причинно-следственную связь между воздействием и конверсией, мы можем начать с основных вопросов и постепенно выстроить статистическую основу, оценивающую причинно-следственные эффекты.

3.9.1.1. Коэффициент конверсии

Один из самых простых вопросов, который можно задать, это измерение простых показателей, таких как *коэффициент конверсии*. Опираясь на предположение, что общее число лиц n , подвергшихся воздействию, известно, и, узнав количество k

конвертировавшихся получателей из числа n , мы можем оценить коэффициент конверсии как

$$R = \frac{k}{n}. \quad (3.68)$$

В зависимости от объема выборки, полученная оценка может быть или не быть статистически достоверной. Если выборка слишком мала, можно ожидать, что полученный коэффициент будет иметь высокую дисперсию и резко изменится, если одну и ту же кампанию провести несколько раз. Если выборка велика, можно ожидать более высокой достоверности результатов. Надежность оценки можно измерить разными способами, с использованием разных статистических методов. В этой книге предпочтение отдается байесовским методам и моделированию Монте-Карло из-за их согласованности и гибкости, поэтому используем этот же подход для рандомизированных экспериментов. Это не самое простое решение для большинства задач, но помогает определить основу, которую можно эффективно распространить на более сложные сценарии, которые мы рассмотрим позже.

Общее количество рекламных акций n — не случайное число, выбранное перед экспериментом, соответственно, наша цель состоит в том, чтобы выяснить распределение коэффициента конверсии с учетом наблюдаемого количества конверсий $p(R | k)$. Если это распределение известно, мы сможем оценить вероятность значительного отклонения результатов от наблюдаемых значений в случае гипотетических повторных экспериментов и тем самым оценить надежность расчетного коэффициента. Согласно правилу Байеса, рассматриваемое распределение можно разложить следующим образом:

$$p(R | k) = \frac{p(k | R)p(R)}{p(k)}, \quad (3.69)$$

где $p(k | R)$ — это правдоподобность, то есть вероятность, наблюдения k конверсий с учетом, что значение коэффициента конверсии известно и равно R , и $p(R)$ — априорное распределение коэффициентов конверсии. Вероятность $p(k)$ можно рассматривать как коэффициент нормализации, поскольку задана точка данных k ; следовательно, этот член просто гарантирует, что распределение коэффициента является распределением вероятности, то есть интеграл по всему диапазону равен 1. Таким образом, это значение можно выразить следующим образом:

$$p(k) = \int p(k | R)p(R)dR. \quad (3.70)$$

Проще говоря, мы начинаем с предварительного предположения о распределении коэффициента $p(R)$, а наблюдаемые данные, то есть количество конвертировав-

шихся клиентов k , подтверждают верность или ложность нашего предположения. Апостериорное распределение $p(R|k)$ получено путем уточнения нашего предположения с учетом наблюдаемых подтверждений.

Поскольку апостериорное распределение коэффициента включает два фактора, $p(k|R)$ и $p(R)$, мы должны задать оба эти распределения. Исходя из предположения о фиксированном значении коэффициента конверсии, вероятность, что ровно k человек из n будут конвертированы, задается биномиальным распределением с функцией вероятности масс вида

$$\begin{aligned} p(k|R) &= \binom{n}{k} \cdot R^k (1-R)^{n-k} = \\ &= \frac{n!}{k!(n-k)!} \cdot R^k (1-R)^{n-k}. \end{aligned} \quad (3.71)$$

Второй фактор, априорное распределение $p(R)$, можно считать однородным или вывести из исторических данных кампании. Рассмотрим сначала случай равномерного распределения. Если априорное распределение $p(R)$ равномерно в диапазоне от 0 до 1, то апостериорное распределение $p(R|k)$ имеет ту же форму, что и вероятность, заданная уравнением 3.71, но теперь является функцией от R , а не от k , поэтому нормализующая константа будет другой. Эту константу можно обозначить как $c(n, k)$ и получить

$$p(R|k) = R^k (1-R)^{n-k} \cdot c(n, k). \quad (3.72)$$

Это распределение известно как *бета-распределение*, и для него существует стандартное обозначение. В этом обозначении апостериор можно выразить как

$$p(R|k) = \text{beta}(k+1, n-k+1), \quad (3.73)$$

где бета-распределение определяется как

$$\begin{aligned} \text{beta}(\alpha, \beta) &= \frac{1}{B(\alpha, \beta)} \cdot x^{\alpha-1} (1-x)^{\beta-1}, \\ B(\alpha, \beta) &= \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx. \end{aligned} \quad (3.74)$$

Распределение коэффициента конверсии по заданному числу человек n , на которых оказано воздействие, и k конверсий описывается бета-распределением.

Если априорное распределение неоднородно, его также можно смоделировать как бета-распределение:

$$p(R) = \text{beta}(x, y), \quad (3.75)$$

где параметры x и y можно определить, например, на основе исторических данных. В этом случае апостериорное распределение по-прежнему является бета-распределением:

$$\begin{aligned} p(R|k) &\propto p(k|R) \cdot p(R) \\ &\propto R^k (1-R)^{n-k} \cdot \text{beta}(x, y) \\ &\propto R^{k+x-1} (1-R)^{n-k+y-1} \\ &\propto \text{beta}(k+x, n-k+y). \end{aligned} \quad (3.76)$$

Говорят, что бета-распределение является *сопряженным априорным распределением* для биномиального распределения: если функция правдоподобия биномиальна, выбор бета-распределения в качестве априорного приведет к тому, что апостериорное распределение также будет бета-распределением. Обратите внимание, что $\text{beta}(1, 1)$ сводится к равномерному распределению, поэтому результат 3.73, полученный для равномерного априорного распределения, является частным случаем выражения 3.76.

Теперь можно оценить вероятность нахождения коэффициента конверсии R в некотором *доверительном интервале* $[a, b]$ как

$$\Pr(a < R < b) = \int_a^b \text{beta}(k+1, n-k+1) dR. \quad (3.77)$$

Уравнение 3.77 можно вычислить аналитически, но точно так же доверительный интервал для коэффициента конверсии можно определить методом Монте-Карло. В этом случае процесс оценки можно описать следующим образом:

1. На входе имеются n , k и желаемый уровень доверия $0 < q < 100$ %.
2. Сгенерировать большое число случайных значений с распределением $\text{beta}(k+1, n-k+1)$.
3. Оценить $q/2$ -й и $(100 - q/2)$ -й процентиля сгенерированных значений, чтобы получить желаемый доверительный интервал. Например, мы можем быть на 95 % уверены, что оценка R лежит между 2,5 % и 97,5 % percentилями.

Примеры бета-распределений для разных значений n и k , а также соответствующие доверительные интервалы показаны на рис. 3.28. Такой подход к моделированию может показаться чрезмерно сложным для оценки основных показателей, таких как коэффициент конверсии, но его преимущества станут более очевидны, когда мы перейдем к более сложным случаям.

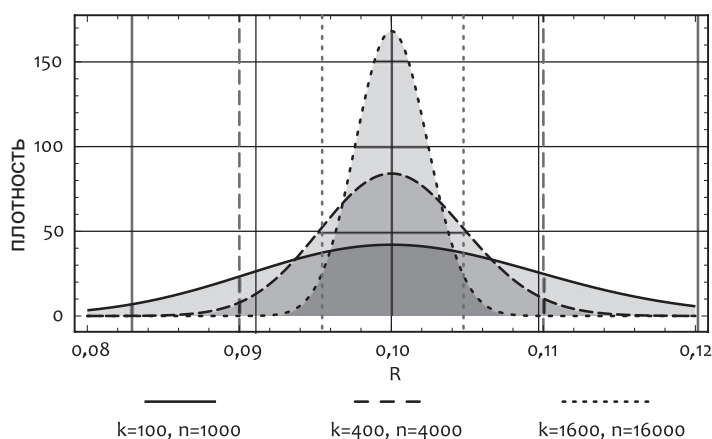


Рис. 3.28. Примеры апостериорного распределения $p(R|k)$ для равномерного априорного распределения и различных размеров выборки n . Среднее $k/n = 0,1$ для всех выборок. Вертикальные линии отмечают 2,5 % и 97,5 % процентили соответствующих распределений. Изначально мы имеем равномерное априорное распределение, и по мере увеличения размера выборки получаем более узкое апостериорное распределение

3.9.1.2. Подъем

Конверсия сама по себе не является достаточным показателем качества алгоритма таргетирования и эффективности маркетинговой кампании. Как уже говорилось выше в этой главе, эффективность обычно измеряется как подъем, который определяется как разность коэффициентов конверсии в опытной и контрольной группах. Коэффициент конверсии в контрольной группе рассматривается как базовый, а повышение можно оценить как коэффициент конверсии в опытной группе, измеренный относительно базового коэффициента:

$$L = \frac{R}{R_0} - 1, \quad (3.78)$$

где R_0 — базовый коэффициент конверсии, а R — коэффициент конверсии для рассматриваемой кампании. Со статистической точки зрения также желательно измерить надежность этой оценки, то есть вероятность

$$\Pr(R > R_0 | data), \quad (3.79)$$

чтобы гарантировать, что полученные результаты обусловлены влиянием рассматриваемой кампании, а не какими-то внешними неконтролируемыми факторами. Стандартный способ решения этой задачи — *рандомизированные эксперименты*. Суть состоит в том, чтобы произвольно разбить потребителей, которые потенци-

ально могут быть вовлечены в кампанию, на две группы (опытную и контрольную), воздействовать на опытную группу (отправлять предложения, показывать рекламу, представить новый дизайн сайта и т. д.), а контрольную группу оставить без воздействия или оказывать на нее обычное базовое воздействие. Случайность отбора потребителей в опытную и контрольную группы играет важную роль, гарантируя исключение влияния на наблюдаемую разницу в результатах систематического смещения между двумя группами, такого как разница в доходах. Для сохранения равенства условий не менее важно наблюдать за опытной и контрольной группами параллельно, чего нельзя обеспечить, например, при сравнении новых данных с историческими.

Схема организации рандомизированных экспериментов для целевых кампаний показана на рис. 3.29. Клиенты с высокой предрасположенностью, выявленные алгоритмом таргетирования, делятся на опытную и контрольную группы, и на опытную группу оказывается воздействие. Количество положительных и отрицательных исходов измеряется для обеих групп: n_T и n_C — количество лиц, k_T и k_C — количество конверсий в опытной и контрольной группах соответственно. Подъем измеряется путем сравнения коэффициента конверсии в опытной группе k_T/n_T с коэффициентом в контрольной группе k_C/n_C .

$$\Pr(a < L < b) = \iint_{a < L < b} L(R_T, R_C) \cdot \Pr(R_T, R_C) dR_T dR_C. \quad (3.80)$$

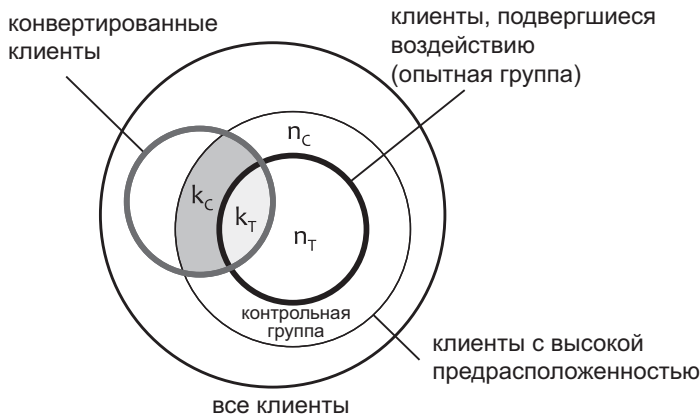


Рис. 3.29. Пример организации рандомизированного эксперимента

Теперь оценим вероятность $\Pr(R_T > R_C)$, то есть найдем доверительный интервал для подъема L . Его можно вычислить так же, как вычислялся доверительный ин-

тервал для коэффициента конверсии в выражении 3.77, но на этот раз мы должны учесть сопряженное распределение для R_T и R_C :

Если допустить, что рандомизированные эксперименты спроектированы и проведены правильно и достигнута независимость между опытной и контрольной группами, можно предположить, что приведенная выше совместная вероятность делится на отдельные распределения коэффициентов конверсии:

$$\Pr(R_T, R_C) = \Pr(R_T | k_T, n_T) \cdot \Pr(R_C | k_C, n_C). \quad (3.81)$$

Теперь можно применить тот же прием имитации, который использовался для отдельного коэффициента конверсии. Коэффициенты R_T и R_C подчиняются бета-распределению, поэтому мы можем сгенерировать образцы подъема, определив два коэффициента конверсии из соответствующих бета-распределений и вычислив отношение. Вот как выглядит этот процесс:

1. Определяются входные значения k_T , n_T , k_C и n_C — по наблюдаемым данным, и желаемый уровень доверия $0 < q < 100$ %.
2. Генерируется большое количество значений L вычислением каждого образца по следующему алгоритму:
 - а) определить R_T из распределения $\text{beta}(k_T + 1, n_T - k_T + 1)$;
 - б) определить R_C из распределения $\text{beta}(k_C + 1, n_C - k_C + 1)$;
 - в) вычислить $L = R_T / R_C - 1$.
3. Оценить желаемый доверительный интервал для L взятием $q/2$ -го и $(100 - q/2)$ -го перцентилей для сгенерированных значений.

Описанный подход с успехом используется во многих практических сценариях — в кампаниях по продвижению, в рекламе и в тестировании произвольных улучшений, таких как новый дизайн веб-сайта. Рандомизированные эксперименты, однако, накладывают определенные ограничения на порядок проведения кампании, и в некоторых случаях это может стать проблемой. В частности, требование к наличию контрольной группы может повлечь дополнительные расходы — мы подробно изучим этот вопрос в следующем разделе.

Важно отметить, что для оценки увеличения доходов не обязательно измерять коэффициент конверсии или даже следить за конверсией отдельных клиентов. Достаточно определить общий доход, полученный в опытной и контрольной группах в течение определенного периода времени после кампании, и оценить подъем как отношение между ними. Это может быть единственным способом оценки подъема, если информация о конверсии недоступна.

3.9.2. Неэкспериментальное исследование

Рандомизированные эксперименты можно использовать в среде онлайн-рекламы для измерения увеличения конверсий, обеспечиваемого кампаниями. Рандомизированные методы требуют осторожности при выборе контрольной группы и гарантируют отсутствие систематического смещения между опытной и контрольной группами. Стандартный подход к достижению этой цели состоит в том, чтобы отложить выбор контрольных пользователей до самого конца конвейера доставки рекламы и отобрать пользователей после этапов таргетирования и торгов, как показано на рис. 3.30. Пользователям из опытной группы открывается доступ к фактическим рекламным объявлениям, а из контрольной группы — к некоторым фиктивным объявлениям, таким как объявления социальной рекламы (Public Service Announcement, PSA), благодаря чему разница между группами может служить мерой влияния рекламы.

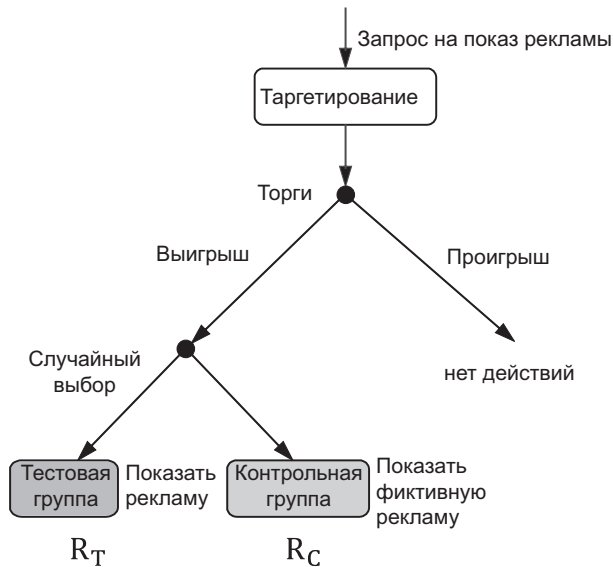


Рис. 3.30. Оценка подъема в онлайн-рекламе с использованием рандомизированных экспериментов. R_T и R_C — это коэффициенты конверсии в опытной и контрольной группах соответственно

Наличие рекламной биржи, однако, представляет серьезную проблему, поскольку показы фиктивных объявлений контрольной группе не производятся бесплатно и должны покупаться, как и фактические показы. Возникает вопрос: можно ли отложить выбор контрольной группы до этапа торгов, как показано на рис. 3.31.

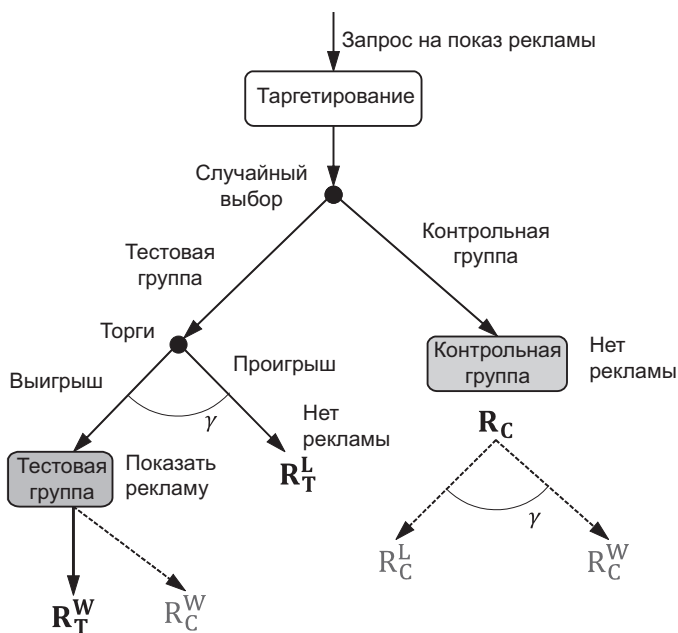


Рис. 3.31. Оценка подъема с использованием методики неэкспериментального исследования

Этот подход фактически означает, что эксперимент перестает быть контролируемым, потому что процесс торгов — какие ставки выиграли, а какие проиграли — не контролируется нами, что, соответственно, может вызвать произвольное смещение в опытной группе по сравнению с контрольной. Мы можем только наблюдать за результатами торгов и конверсиями и измерять эффект рекламы путем статистических выводов. Это приводит нас в обширную область теории неэкспериментальных исследований и выводов о причинной зависимости, которая интенсивно развивалась на протяжении десятилетий и обусловлена необходимостью анализа процессов, неподконтрольных исследователям. Наша проблема с предвзятостью торгов близко соответствует проблеме *эффекта лечения при несоответствии* в клинических испытаниях. Эффект лечения можно оценить с помощью рандомизированных экспериментов и сравнения испытуемых из опытной группы, подвергшихся лечению, с испытуемыми из контрольной группы. Даже при том, что испытуемые могут произвольно включаться в опытную и контрольную группы, некоторые члены опытной группы не могут подвергаться лечению из-за проблемы несоответствия. Разделение на совместимые и несовместимые подгруппы после рандомизации соответствует разделению выигрыша и проигрыша в процессе торгов, когда оно следует за выбором контрольной группы, поэтому мы можем

использовать методологию исследований, предназначенную для клинических испытаний с несоответствием.

Задачу оценки подъема с помощью неэкспериментальных исследований можно решить с помощью различных методов. Начнем с базового метода, который иллюстрирует, как некоторые идеи теории причинности можно применить к задаче [Chalasani and Sriharsha, 2016; Rubin, 1974; Jo, 2002].

На рис. 3.31 можно видеть, что у нас есть по меньшей мере три коэффициента конверсии, которые можно измерить непосредственно: R_C для контрольной группы, R_T^L для проигравших ставок в опытной группе и R_T^W для пользователей, подвергшихся фактическому воздействию. Наша цель — найти коэффициент R_C^W , который можно интерпретировать как коэффициент *потенциальной* конверсии пользователей, которые конвертировались бы даже без воздействия. Это гипотетическое значение, потому что мы не можем уйти в прошлое, отменить воздействие и посмотреть, что из этого получится. Однако его можно оценить по известным данным при определенных допущениях. Во-первых, отметим, что соотношение γ между количеством «выигравших» и «проигравших» пользователей можно наблюдать непосредственно. Предположив, что распределение тех и других одинаково в опытной и в контрольной группах, можно утверждать, что

$$R_C = \gamma \cdot R_C^W + (1 - \gamma) R_C^L, \quad (3.82)$$

где R_C^W и R_C^L являются коэффициентами конверсии для пользователей в контрольной группе, которые могут выиграть и проиграть соответственно, если их отобрать в опытную группу. Второе предположение состоит в том, что $R_C^L = R_T^L$, потому что обе группы содержат только «проигравших» пользователей, не подвергавшихся воздействию рекламы, соответственно мы не ожидаем никакого смещения между ними. Следовательно, R_C^W можно выразить, используя известные значения:

$$R_C^W = \frac{1}{\gamma} \cdot (R_C - (1 - \gamma) R_T^L). \quad (3.83)$$

Наконец, подъем можно оценить как отношение между наблюдаемым R_T^W и выведенным R_C^W .

Надежность оценки подъема можно оценить с помощью того же метода имитации, который использовался для рандомизированных экспериментов. Для этого требуется сгенерировать выборки в соответствии с распределением подъема, которое иногда бывает сложно задать, потому что обусловлено несколькими случайными процессами: выбор контрольной группы, торги и конверсия. Мы наблюдаем только часть информации для каждой реализации этого сложного процесса (выбранная группа, результат торгов и результат конверсии), но мы не наблюдаем внутренних

свойств пользователей и других *скрытых факторов*, которые определяют сопряженное распределение вероятностей наблюдаемых результатов. В остальной части этого раздела мы обсудим статистический аппарат, сочетающий идею потенциальных результатов, рассмотренных выше, с продвинутыми методами имитации для вывода распределений различных свойств кампании, включая подъем [Chickering and Pearl, 1996]. Мы опишем этот аппарат в два этапа. Сначала уточним модель интересующих нас случайных процессов. А затем посмотрим, как оценить модель с помощью имитации.

3.9.2.1. Описание модели

Объяснить скрытые факторы и их влияние можно с помощью графической модели, представленной на рис. 3.32. Каждый узел представляет случайную переменную, а стрелки соответствуют зависимостям между узлами. Случайные переменные Z , A и Y соответствуют рандомизации, торгам и конверсии. Точнее, бинарная переменная $Z \in \{0, 1\}$ принимает значение 1, если пользователь выбран в контрольную группу, и 0 в противном случае. Переменная $A \in \{0, 1\}$ принимает значение 1, если мы выиграли торги и показали объявление, и 0 в противном случае. И, наконец, переменная $Y \in \{0, 1\}$ принимает значение 1, если пользователь конвертировался, и 0 в противном случае. Случайная переменная S соответствует состоянию пользователя и, возможно, другим скрытым факторам, влияющим на способность рекламодателя выиграть ставку и получить отклик после показа.

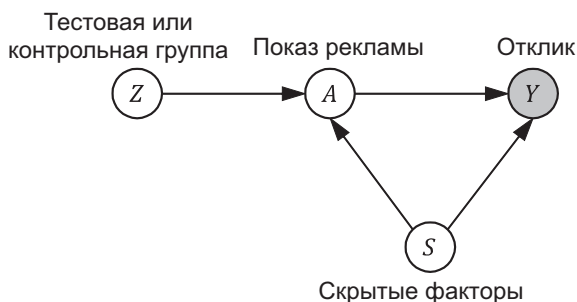


Рис. 3.32. Графическая модель неэкспериментального исследования со скрытыми факторами

Прежде всего, мы должны выяснить параметры сопряженного распределения $\text{Pr}(z, a, y, s)$ и проинтегрировать его, чтобы получить доверительный интервал для подъема L . Вопрос в том, как разложить распределение $\text{Pr}(z, a, y, s)$, чтобы его можно было обрабатывать на компьютере. Графическая модель на рис. 3.32 делает определенные предположения, которые можно использовать для декомпозиции:

Z и S считаются независимыми, так как рандомизация не должна зависеть от внешних факторов, а Z и Y условно независимы от A и S , потому что на конверсию можно повлиять только через события A . Это приводит к следующему разложению плотности вероятности:

$$\Pr(z, a, y, s) = \Pr(z)\Pr(a)\Pr(a|z, s)\Pr(y|a, s). \quad (3.84)$$

Теперь необходимо определить случайную переменную состояния S и ее роль в плотностях $\Pr(a|z, s)$ и $\Pr(y|a, s)$. Идея скрытых факторов состоит в том, чтобы учесть «состояние мира», которое не наблюдается непосредственно, но может влиять на результаты, такие как подъем. Эту идею можно рассматривать как аналог потенциальных результатов, которые мы рассматривали в начале этого раздела, потому что, сумев вывести состояние из наблюдений, мы сумеем оценить потенциальные результаты для разных предварительных условий. Например, если известно, что данный пользователь никогда не выиграет на бирже, мы сможем предсказать результаты включения этого пользователя в опытную или в контрольную группу.

Скрытое состояние можно смоделировать по-разному, в зависимости от доступных данных, интересующих нас показателей и общего понимания предметной области. Мы используем стандартную модель, иллюстрирующую, как скрытые состояния можно определить в виде функции наблюдаемых данных и как показатели, такие как подъем, можно получить из состояний [Heckerman and Shachter, 1995; Chickering and Pearl, 1996].

С точки зрения эффективности кампании нас интересуют в основном два свойства пользователя: соответствие рекламному методу (возможность или невозможность выиграть ставку) и отклик на рекламу (конвертируется или нет). Эти свойства соответствуют описанным выше вероятностям $\Pr(a|z, s)$ и $\Pr(y|a, s)$ и могут рассматриваться как внутреннее состояние пользователя, систематически влияющее на полученные результаты. Мы можем перечислить возможные состояния отдельно для соответствия и отклика и задать условие для каждого состояния, указывающее на возможность данного состояния для наблюдаемого кортежа (z, a, y) .

Набор возможных состояний пользователя является декартовым произведением поведений соответствия и отклика, что дает нам множество с 16 элементами $\{s_1, \dots, s_{16}\}$, представляющими все пары (C_p, R_q) поведений соответствия и отклика, перечисленные в табл. 3.10 и 3.11:

$$\begin{aligned} S &\in \{s_1, \dots, s_{16}\}, \\ s_{p+4(q-1)} &= (C_p, R_q), \quad 1 \leq p, q \leq 4. \end{aligned} \quad (3.85)$$

Таблица 3.10. Состояния соответствий пользователя и условия. Состояния C_3 и C_4 не должны иметь места в рассматриваемом сценарии, но они могут возникать в других средах, таких как многоканальная реклама

Соответствие	Условие	Описание
C_1	$a = 0$	Пользователь никогда не видел рекламы
C_2	$a = z$	Пользователь видит рекламу всякий раз, когда наша ставка выигрывает, и только тогда
C_3	$a \neq z$	Пользователь видит рекламу, только если наша ставка не выиграла
C_4	$a = 1$	Пользователь всегда видит рекламу

Таблица 3.11. Состояния откликов пользователя и условия

Отклик	Условие	Описание
R_1	$y = 0$	Пользователь никогда не конвертируется
R_2	$y = a$	Пользователь конвертируется только после воздействия
R_3	$y \neq a$	Пользователь конвертируется, только если не было воздействия
R_4	$y = 1$	Пользователь всегда конвертируется

Следовательно, S — это случайная переменная с 16 состояниями, взятая из множества 16 возможных состояний.

Мы непосредственно наблюдаем бинарные кортежи (z^j, a^j, y^j) для каждого пользователя j , но пользовательские состояния s^j никогда не наблюдаются непосредственно. Однако, если вывести состояние, на его основе можно оценить интересующие нас потенциальные параметры, такие как подъем. В частности, нас интересуют не отдельные пользовательские состояния, а вектор долей состояний:

$$\mu = (\mu_1, \dots, \mu_{16}), \quad (3.86)$$

где каждая доля μ_i — это отношение числа пользователей в соответствующем состоянии s_i к общему числу наблюдаемых пользователей. Теперь показатели можно определить как функции от μ . Например, подъем $L(\mu)$ можно определить как отношение суммы четырех значений μ , соответствующие состояниям с компонентом отклика R_2 (и любым компонентом соответствия) к сумме других четырех значений μ , соответствующих состояниям с компонентом отклика R_3 . Для ответов на другие вопросы можно также определить другие функции от μ .

3.9.2.2. Имитация

Взяв за основу модель, описанную выше, доверительный интервал показателя $L(\mu)$ можно выразить через апостериорное распределение случайного вектора μ :

$$\Pr(a < L(\mu) < b) = \int_{a < L(\mu) < b} L(\mu) \cdot p(\mu | \text{data}) d\mu, \quad (3.87)$$

где data представляет все наблюдаемые кортежи (z^j, a^j, y^j) . Обозначим вектор состояний пользователей как

$$\mathbf{s} = (s^1, \dots, s^m), \quad (3.88)$$

где m — число наблюдаемых пользователей. В таком случае распределение долей состояний μ можно рассматривать как случайную функцию состояний пользователей \mathbf{s} , которые, в свою очередь, также являются случайными переменными — они не наблюдаются явно, но их вероятностные характеристики можно вывести из данных. Следовательно, мы должны рассмотреть сопряженное распределение 16 переменных в переменных μ и m из \mathbf{s} :

$$\Pr(a < L(\mu) < b) = \int_{a < L(\mu) < b} L(\mu) \cdot p(\mu, \mathbf{s} | \text{data}) d\mu d\mathbf{s}. \quad (3.89)$$

Метод имитации требует аппроксимации распределения $p(\mu, \mathbf{s} | \text{data})$ на основе наблюдаемых данных, то есть должна иметься возможность получить векторы μ из этого распределения. После получения векторов можно вычислить образцы $L(\mu)$ и оценить их распределение. Теперь мы должны ответить на вопрос: как получить образцы из распределения $p(\mu, \mathbf{s} | \text{data})$. Мы не знаем функциональной формы распределения, но существуют статистические методы, способные помочь нам в создании экземпляров из распределения, не определяя его явно.

Для извлечения образцов из многомерных распределений широко используется метод семплирования по Гиббсу¹ (Gibbs sampling) [Geman and Geman, 1984]. Предположим, что нам нужно взять образцы из многомерного распределения $p(x_1, \dots, x_n)$. В методе Гиббса используется тот факт, что многомерное распределение можно разделить на n условных распределений:

$$p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n), \quad 1 \leq i \leq n. \quad (3.90)$$

¹ https://ru.wikipedia.org/wiki/Семплирование_по_Гиббсу

Иногда невозможно отобрать точки непосредственно из многомерного распределения, но выборка из условного распределения возможна всегда. Идея метода семплирования Гиббса состоит в том, чтобы вместо вероятностного выбора всех n переменных одновременно выбирать одну переменную за раз, зафиксировав текущие значения в оставшихся переменных. Другими словами, каждая переменная выбирается из ее условного распределения с остальными фиксированными переменными:

$$x_j \sim p(x_j | x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n), \quad 1 \leq j \leq n. \quad (3.91)$$

Этот итеративный алгоритм многократно извлекает образцы из условных распределений, подставляя в условия образцы, сгенерированные ранее. Например, рассмотрим базовый случай с двумя переменными x_1 и x_2 . Переменные сначала инициализируются некоторыми значениями, которые можно получить из априорного распределения, а затем обновляются в каждой итерации i в соответствии со следующими правилами:

$$\begin{aligned} x_1^{(i)} &\sim p(x_1 | x_2^{(i-1)}), \\ x_2^{(i)} &\sim p(x_2 | x_1^{(i)}). \end{aligned} \quad (3.92)$$

Для сходимости этому процессу может потребоваться определенное количество итераций, после чего он начнет создавать точки, подчиняющиеся распределению $p(x_1, x_2)$. Этот метод очень эффективен на практике, потому что условные распределения часто гораздо проще определить, чем сопряженные. Обобщенная версия семплера Гиббса представлена в алгоритме 3.1.

Алгоритм 3.1. Семплер Гиббса

инициализировать $(x_1^{(0)}, \dots, x_n^{(0)})$ из априорного распределения

для итерации $i = 1, 2, \dots$ **do**

	выбрать $x_1^{(i)} \sim p(x_1 x_2^{(i-1)}, x_3^{(i-1)}, \dots, x_n^{(i-1)})$
	выбрать $x_2^{(i)} \sim p(x_2 x_1^{(i)}, x_3^{(i-1)}, \dots, x_n^{(i-1)})$
	...
	выбрать $x_n^{(i)} \sim p(x_n x_1^{(i)}, x_2^{(i)}, \dots, x_{n-1}^{(i)})$

конец

Вернемся теперь к распределению $p(\mu, s | data)$ и посмотрим, как выбрать образцы из него с помощью семплера Гиббса.

Поскольку семплер отдельно выбирает каждый элемент μ и \mathbf{s} , можно задать отдельные процедуры оценки для $p(\mathbf{s} | \mu, data)$ и $p(\mu | \mathbf{s}, data)$.

Для первой вероятности можно использовать предположение о независимости пользователей, поэтому апостериорные вероятности их состояний задаются как

$$p(s^j = s_i | \mu, \mathbf{s}, data) \propto p(a^j, y^j | z^j, s_i) \cdot \mu_i, \quad (3.93)$$

где $p(a^j, y^j | z^j, s_i)$ — вероятность наблюдения исходов z^j , a^j и y^j с учетом состояния s_i . Можно предположить, что вероятность равна единице, если наблюдаемые исходы согласуются с условиями состояния s_i , в противном случае она равна нулю. Следовательно, вероятность состояния s_i для пользователя j можно оценить на основе известных значений a^j , y^j и z^j и условий состояния из табл. 3.10 и 3.11. В модели с 16 состояниями для каждого пользователя требуется получить вектор с 16 вероятностями. Затем этот вектор умножается на априорную вероятность состояния i в соответствии с правой частью выражения 3.93. Получившийся вектор с 16 числами определяет полиномиальное распределение, из которого можно извлечь образец s^j .

Вторая часть — это условное распределение $p(\mu | \mathbf{s}, data)$. Обозначим через n_i , сколько раз состояние s_i встречается в \mathbf{s} . Поскольку μ является вектором долей состояний, то есть каждый элемент i является эмпирической вероятностью состояния s_i , вектор счетчиков n_i будет иметь полиномиальное распределение с параметром μ . Следовательно, вероятность наблюдения вектора \mathbf{s} для данного вектора долей состояний μ равна

$$\prod_i \mu_i^{n_i}, \quad (3.94)$$

откуда апостериорное распределение долей состояний будет определяться как

$$p(\mu | \mathbf{s}, data) \propto \prod_i \mu_i^{n_i} \cdot \Pr(\mu). \quad (3.95)$$

Последний шаг — задать априорное распределение $\Pr(\mu)$. Выше в рандомизированных экспериментах мы использовали бета-распределение, потому что вероятность имела биномиальное распределение, а сопряженным априорным распределением для биномиального было бета-распределение. Аналогично мы имеем теперь полиномиальное распределение, сопряженным априором для которого является распределение Дирихле (см. приложение А): если для аппроксимации $\Pr(\mu)$ выбрать распределение Дирихле, апостериорным распределением в выражении 3.96 также будет распределение Дирихле. Более формально предпосылочное убеждение можно выразить как набор счетчиков n_i^0 , которые используются в качестве параметров априорного распределения Дирихле, а апостериорное можно выразить как

$$\begin{aligned}
p(\boldsymbol{\mu} | \mathbf{s}, \text{data}) &\propto \prod_i \mu_j^{n_i} \cdot \text{Dir}(n_1^0, \dots, n_{16}^0) \\
&\propto \prod_i \mu_j^{n_i^0 + n_i - 1} \\
&\propto \text{Dir}(n_1^0 + n_1, \dots, n_{16}^0 + n_{16}).
\end{aligned} \tag{3.96}$$

Уравнения выше можно включить непосредственно к семплер Гиббса: сгенерировать образцы $\boldsymbol{\mu}$ с помощью выражения 3.96, сгенерировать m образцов \mathbf{s} с помощью уравнения 3.93, а затем повторить процесс итеративно, пока не получим достаточное количество реализаций вектора $\boldsymbol{\mu}$ для оценки доверительного интервала $L(\boldsymbol{\mu})$.

3.10. Архитектура систем таргетирования

Системы таргетирования можно реализовать по-разному, в зависимости от конкретной сферы и применения. Однако некоторые логические компоненты остаются общими для большинства систем таргетирования. В этом разделе мы рассмотрим каноническую архитектуру, включающую все основные логические блоки, необходимые для создания целевой рекламы или рекламных акций. Эта архитектура предполагает, что система функционирует в режиме «запрос/ответ», то есть в реальном времени получает запросы, содержащие некоторую контекстную информацию, такую как идентификатор потребителя или идентификатор канала, и возвращает одно или несколько предложений для этого конкретного контекста. Мы считаем это приложение и его дизайн наиболее универсальными и важными, однако его можно адаптировать к другим приложениям, таким как пакетная рассылка электронной почты.

На рис. 3.33 показана обобщенная логическая архитектура системы таргетирования. Она предполагает наличие трех основных подсистем, каждая из которых состоит из нескольких компонентов.

3.10.1. Сервер таргетирования

Сервер таргетирования инкапсулирует основную логику обработки входящих запросов и выдачи ответов с рекламными предложениями. Его можно рассматривать как конвейер со следующими основными этапами:

УСЛОВИЯ. В соответствии с описанным выше процессом таргетирования на первом этапе выполняется проверка явных ограничений для всех возможных рекламных предложений. Примерами таких условий могут служить наличие определенных товаров в текущей корзине, определенные географические координаты и т. д.

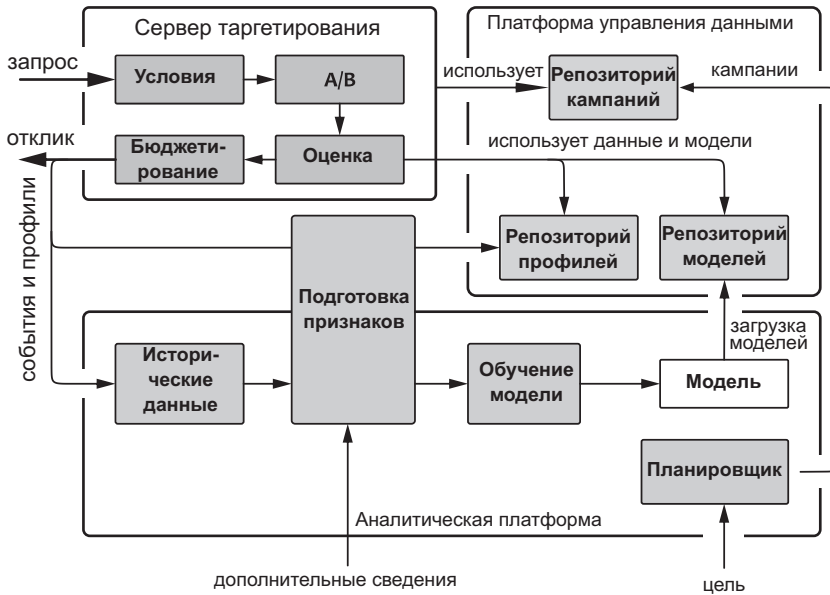


Рис. 3.33. Обобщенная архитектура системы таргетирования

А/В ТЕСТИРОВАНИЕ. Одновременная оценка сразу нескольких стратегий таргетирования является стандартным методом, который помогает бороться с неточностью данных и моделей, а также измерять эффективность кампаний. Сервер таргетирования может назначать разные методы рекламы разным пользователям и сообщать показатели эффективности для каждого метода отдельно, благодаря чему оптимальную стратегию можно выбрать позже. А/В тестирование может проводиться для разных аспектов процесса таргетирования и потребительского опыта, таких как разные методы оценки, разные текстовые сообщения или изображения и т. д.

Выбор стратегии, как правило, производится для каждого потребителя. После выбора стратегия сохраняется в профиле потребителя и последовательно используется для всех запросов, связанных с ним. Это помогает добиться непротиворечивого пользовательского опыта и создать для разных стратегий непересекающиеся потребительские сегменты.

Также нередко одна или несколько экспериментальных стратегий сравнивается с контрольной группой — группой потребителей, получающих базовое воздействие, которое можно рассматривать как базовую стратегию или ее отсутствие. Эффективность экспериментальных стратегий затем можно оценить как подъем относительно базовой линии.

ОЦЕНКА. Сервер таргетирования оценивает стимулы, которые прошли предыдущий этап, путем оценки моделей предрасположенности, связанных контекстом со стимулом, включая исторический профиль потребителя. Модели предрасположенности можно динамически выбирать под стимул на основе метаданных и бизнес-правил, чтобы избежать ручной привязки модели к каждому стимулу.

ПРАВИЛА БЮДЖЕТИРОВАНИЯ. Окончательный ответ создается на основе списка оценок допустимых поощрений путем применения правил бюджетирования и других ограничений на количество показов для данного потребителя, канала или кампании.

3.10.2. Платформа управления данными

Платформу управления данными можно рассматривать как рабочую базу данных для хранения профилей клиентов и других данных, необходимых для таргетирования, включая конфигурации кампаний. Основными компонентами платформы управления данными можно считать:

РЕПОЗИТОРИЙ ПРОФИЛЕЙ. Репозиторий, в котором хранятся исторические данные об отдельных потребителях, в том числе исходные данные, такие как отдельные заказы или события на веб-сайте, и агрегированные статистики (признаки), которые можно непосредственно использовать в моделях оценки предрасположенности. Репозиторий может быть заполнен контекстными данными с сервера таргетирования и из внешних источников данных.

РЕПОЗИТОРИЙ МОДЕЛЕЙ. Репозиторий, в котором хранится пул моделей предрасположенности, используемых сервером таргетирования. Эти модели создаются и обновляются аналитической платформой, описанной ниже.

РЕПОЗИТОРИЙ КАМПАНИЙ. Репозиторий, в котором хранятся сведения о конфигурации кампаний, включая графические ресурсы, условия, ограничения бюджета и т. д.

3.10.3. Аналитическая платформа

Аналитическая платформа собирает, объединяет и сохраняет данные профиля клиента, а также дополнительные сведения, необходимые для моделирования, подготовки данных и отчетности. Примерами таких сведений могут служить: информация из каталога товаров, данные о продажах и магазине и т. д. Одним из основных результатов процесса подготовки данных являются *признаки профиля*, которые можно использовать для обучения и оценки предиктивных моделей. Данные, объединенные и подготовленные платформой, могут предоставляться разным

потребителям с разными соглашениями об уровне обслуживания. Например, для аналитических целей признаки профилей можно создавать в пакетном режиме, но для оценки модели те же признаки должны передаваться серверу таргетирования в реальном времени. То есть некоторые модули подготовки данных могут предусматривать работу в разных режимах. Этот аспект иллюстрируется блоком подготовки признаков на рис. 3.33, который используется и для моделирования, и для агрегирования данных в реальном времени с последующей передачей платформе управления данными.

Одной из основных функций платформы является создание моделей предрасположенности путем запуска алгоритмов машинного обучения с данными, поступающими из сервера таргетирования, и внешних источников данных. Платформа может включать инструменты для создания моделей вручную и их автоматического обновления (переобучения). Кроме того, аналитическая система выполняет измерения и дает возможность создания отчетов и исследовательского анализа данных.

Наконец, аналитическая платформа может содержать планировщик. Это ключевой компонент программной маркетинговой системы, проектирующий и оптимизирующий рекламные кампании. Планировщик использует исторические данные и статистику, бизнес-правила, передовые методики и эвристики и определяет оптимальные стратегии (продолжительность и тип стимулирования, каналы, модели предрасположенности и т. д.) на основе цели и дополнительных ограничений, таких как лимиты бюджета. Он также может спрогнозировать эффективность кампании на основе исторических данных. Планировщик может иметь несколько функциональных блоков.

ИНВЕСТИЦИОННЫЙ ПЛАНИРОВЩИК. Инвестиционный планировщик получает обобщенное представление о возможностях рынка из исторических данных. Это помогает конечному пользователю правильно определить бизнес-цели и распределить бюджет по разным стратегическим направлениям и кампаниям. Его можно рассматривать как глобальный инструмент оптимизации.

ПЛАНИРОВЩИК КАМПАНИЙ. Планировщик кампаний оптимизирует отдельные рекламные кампании и акции, предлагаемые инвестиционным планировщиком. Вычисляет оптимальные сроки, расходы и т. д.

3.11. Итоги

- Услуги продвижения и рекламы основное внимание уделяют проблеме таргетирования, то есть поиску оптимального соответствия между потребителями и предложениями. В зависимости от приложения системе может потребовать-

ся отобрать предложения для данного потребителя или потребителей для данного предложения.

- Услуги продвижения и рекламы, как правило, руководствуются целью оптимизировать доходы, но важно также помнить об оптимизации качества обслуживания клиентов.
- Основными бизнес-средами для услуг продвижения и рекламы являются стимулирование сбыта и онлайн-реклама. Основными действующими лицами и субъектами в среде продаж являются потребители, производители, ретейлеры и маркетинговые кампании.
- Бизнес-цели производителей и ретейлеров можно смоделировать в терминах затрат и выгод кампании. К непосредственным целям относится, например, прибыльность кампании, а стратегические цели можно описать в терминах жизненного цикла клиента. Ключевыми стратегическими целями являются: привлечение новых клиентов, а также максимизация и удержание существующих клиентов.
- Система таргетирования может иметь вид конвейера, который начинается с распределения ресурсов между целями с последующей подгонкой шаблонов кампаний к цели, связыванием моделей таргетирования и, наконец, выполнением кампаний.
- Платформа моделирования откликов учитывает затраты на кампанию, доходы и статистические свойства клиентов. Основной принцип моделирования отклика — максимизация подъема, то есть прироста прибыли. Подъем можно измерить постфактум, путем использования опытных и контрольных групп.
- Роль основных строительных блоков в системах таргетирования играют модели предрасположенности, времени до наступления события и пожизненной ценности. Основными примерами моделей таргетирования являются уровни лояльности и RFM-анализ. Эти подходы ориентированы на финансовые результаты и не исследуют глубоко причины того или иного поведения клиентов.
- Целью моделирования предрасположенности является поиск потребителей, которые с относительно высокой вероятностью будут вести себя определенным образом, например купят новый продукт. Подобное моделирование является одним из основных методов моделирования предрасположенности.
- Часто удобнее измерить время до наступления события, чем вероятность этого события. Это можно сделать с помощью анализа выживаемости. Модели выживаемости, как и регрессионные модели, могут выражать время до события как функцию независимых переменных, таких как свойства клиента или величина скидки.

- Модели пожизненной ценности оценивают общую сумму денег, которую бренд, вероятно, получит от данного клиента в течение срока их отношений. LTV-моделирование возможно с использованием описательных и предиктивных подходов.
- Маркетинговую кампанию можно сконструировать из нескольких строительных блоков. Шаблон кампании может включать условия таргетирования, модели оценки, правила бюджетирования и ограничения. Кампания обычно соответствует определенной точке или набору точек в путешествии клиента и имеет целью повлиять на него. В качестве примеров шаблонов кампаний можно привести кампании по продвижению продуктов, многоступенчатые кампании, а также кампании удержания и пополнения.
- Таргетирование можно рассматривать как распределение ресурсов между клиентами, но их также можно распределять по каналам, целям, территориям и другим критериям.
- Многие методы и принципы таргетирования продвижения применимы также в других областях, например в онлайн-рекламе, но каждая область имеет свои бизнес-цели и сложности в их реализации. Цели онлайн-рекламы часто определяются в терминах затрат на действие и оценки эффективности конкурирующих рекламодателей.
- Онлайн-реклама использует большое количество приемов и методов таргетирования. Многие из них основаны на моделировании подобию, понятии близости бренда, вероятности отклика и качества запасов.
- Эффективность продвижения и рекламных кампаний обычно измеряется с помощью рандомизированных экспериментов, с созданием опытных и контрольных групп. В некоторых средах, включая онлайн-рекламу, организация контрольных групп связана с дополнительными затратами или упущенной выгодой, поэтому эффективность можно измерить с помощью более совершенных методов неэкспериментальных исследований.

4

Поиск

Таргетирование и реклама решают задачу поиска правильной аудитории для данного продукта или услуги. Не менее важным аналогом таргетирования является *открытие продукта*¹ (product discovery), то есть задача предоставления клиентам удобных механизмов и интерфейсов для просмотра ассортимента и поиска нужных им продуктов. Сервисы таргетирования и открытия — это два основных инструмента программного маркетинга, которые можно использовать для улучшения осведомленности клиентов о продукте, услуге или бренде.

Проблема открытия продукта вращается вокруг понятия *намерения совершения покупки* (purchasing intent). Иногда клиенты явно выражают свое намерение, вводя поисковый запрос или указывая нужные атрибуты продукта каким-то другим способом. Иногда намерение не выражено явно, и программная служба должна определить его по известным характеристикам клиента и его поведению. Эти два сценария часто различаются и рассматриваются как две разные категории служб — службы поиска помогают клиентам находить продукты, отвечающие явно выраженным требованиям, в то время как службы рекомендаций не требуют от пользователей явно выражать свои намерения. Однако граница между поиском и рекомендациями весьма условна. Простая служба поиска может использовать для поиска продуктов только явно введенный запрос. Более продвинутое решение может использовать дополнительные сведения о пользователе, чтобы персонализировать поиск. В некоторых приложениях такие неявные сигналы могут стать важнее явных, и служба поиска структурно трансформируется в службу рекомендаций. С этой точки зрения службы поиска и рекомендации можно сравнить с услугами продавца в магазине, который может найти конкретный продукт по запросу или предложить варианты, исходя лишь из информации о клиенте и его потребностях.

¹ Также может называться «исследование продукта». Целью этого этапа является разработка концепции продукта. — *Примеч. ред.*

Поиск и рекомендации важны не только как функциональные услуги, но и как фундаментальная возможность объединить множество сильных, слабых или посторонних сигналов, чтобы правильно понять потребности клиента и определить продукты, соответствующие этим потребностям. Эта возможность является ключом к созданию эффективных клиентских служб и приложений.

Первый сервис открытия продуктов, который мы рассмотрим, — это поиск. Целью поиска является выбор предложений, соответствующих намерениям клиента, выраженным в поисковом запросе или определяемым выбранными фильтрами. Задачи этого типа решаются с помощью теории поиска информации, поэтому в нашем распоряжении имеется широкий спектр теоретических основ и практических методов поиска. Основная цель этого раздела — отобрать и адаптировать приемы и методы, относящиеся к маркетинговым приложениям. Некоторые из этих методов заимствованы из универсального инструментария, созданного в рамках теории информационного поиска, другие были разработаны специально для маркетинга и продвижения. Мы применим практический подход к методам поиска и сосредоточимся на опыте практического применения, методах и примерах, а не на теории информации. В то же время постараемся по возможности избегать деталей реализации, таких как индексация данных, и сосредоточимся на коммерческой ценности, обеспечиваемой результатами релевантного поиска.

Начнем эту главу с обзора среды и экономических целей. Затем покажем, что задачу релевантного поиска можно выразить в терминах признаков, сигналов и элементов управления, подобно другим программным службам. Попутно рассмотрим методы проектирования, смешивания и настройки этих сигналов и элементов управления в ручном режиме, а затем обсудим, как использовать предиктивную аналитику для автоматической оптимизации.

4.1. Среда

Поиск — один из самых удобных и естественных интерфейсов между человеком и компьютером. Поиск является неотъемлемой частью широкого спектра сервисов и приложений в различных областях, использующих функции поиска по-разному. Эти среды могут существенно различаться как экономическими целями, так и техническими свойствами, такими как объем данных, и могут существенно влиять на разработку и реализацию служб поиска. Рассмотрим несколько основных примеров.

ВЕБ-ПОИСК. Всемирная паутина, или веб, — это набор веб-страниц, содержащих текстовую и мультимедийную информацию, поэтому неудивительно, что веб-поиск изначально рассматривался как задача анализа текста. Первые

поисковые системы главным образом основывались на оценке содержимого страниц, что давало владельцам сайтов широкие возможности обманывать системы поиска, манипулируя скрытыми полями и используя другие методы повышения релевантности. Создатели поисковых систем были вынуждены разработать совершенно новую стратегию, использующую перекрестные ссылки между веб-сайтами и уровень доверия к ним как основной оценочный сигнал, поэтому веб-домены и страницы, на которые ссылаются надежные и популярные сайты, оцениваются высоко, а релевантность ресурсов без ссылок резко снижается. Этот подход, впервые предложенный компанией Google в виде алгоритма PageRank, стал важной чертой веб-поиска, отличающей его от большинства других поисковых сред. Большие масштабы Всемирной паутины и необходимость индексирования огромных объемов данных также значительно повлияли на методы веб-поиска.

ПОИСК ДЛЯ ПРОДВИЖЕНИЯ. Во многих поисковых приложениях организационная релевантность дополняется или переопределяется бизнес-правилами, направленными на достижение бизнес-целей и соблюдение ограничений. Ярким примером таких сред является поиск для продвижения в электронной коммерции, розничной торговле и других потребительских приложениях, таких как службы бронирования гостиниц или поиска ресторанов. Эти приложения требуют функции поиска, которые могут увеличить прибыль за счет рекламы высокодоходных продуктов, помощи в распродаже товаров с истекающим сроком действия или продвижения спонсируемых предложений. Поиск для продвижения должен также учитывать специфическую терминологию и шаблоны использования, чтобы распознать намерения пользователя и общие идиомы.

ЭКСПЕРТНЫЙ ПОИСК. Большое количество поисковых служб и приложений в области права, медицины, научных исследований и промышленности относится к категории экспертного поиска. Экспертный поиск, используемый профессионалами, требует глубокого понимания предметной области, включая жаргон и скрытые отношения между понятиями. Кроме того, экспертный поиск должен поддерживать специализированные шаблоны использования и определения релевантности, которые могут сильно отличаться от типичных потребностей в веб-поиске и в поиске для продвижения товаров. Например, юристу или специалисту по патентам может быть важно найти и изучить каждый документ, относящийся к рассматриваемой теме, тогда как пользователям веб-поиска или поиска для продвижения товаров обычно достаточно нескольких наиболее релевантных результатов. Иногда экспертный поиск также называют *корпоративным поиском*.

Несмотря на все различия, поисковые решения для разных областей имеют много общего. В этой главе мы сосредоточимся на поиске для продвижения товаров, потому что он лучше всего соответствует замыслу этой книги, хотя большинство

строительных блоков, которые будут описаны ниже, вполне применимы в других областях.

Начнем с описания минимальной среды для поиска продвижения, как показано на рис. 4.1, которая вводит ключевые сущности и допущения:

- Основное назначение службы поиска для продвижения — дать клиентам простой интерфейс для поиска товаров, услуг и прочего с использованием запросов в свободной текстовой форме и дополнительной контекстной информации, которая может включать профиль клиента, тип устройства просмотра, географическое местоположение и т. д. Другими словами, входные данные службы поиска, как предполагается, состоят из пары запрос/контекст, а выходные данные представляют список сущностей в порядке уменьшения релевантности. Эта базовая функциональность может быть расширена дополнительными инструментами поиска, которые мы подробно рассмотрим позже.
- Все результаты поиска обычно представлены в виде сочетания структурированных или полуструктурированных записей из нескольких источников. Например, онлайн-магазин может извлекать номенклатуру товаров из текстовых описаний, предоставленных производителями, оценок и отзывов пользователей, данных о продажах, данных о запасах, местоположений магазинов, прейскурантов и внутренних метаданных, таких как иерархии категорий.
- Механизм поиска выдает результаты, сопоставляя запрос или характеристики контекста с признаками сущностей. Этот процесс можно настроить с помощью нескольких средств управления релевантностью, которые определяют порядок извлечения признаков из исходных данных, их сопоставления и смешивания различных сигналов для получения окончательных результатов.

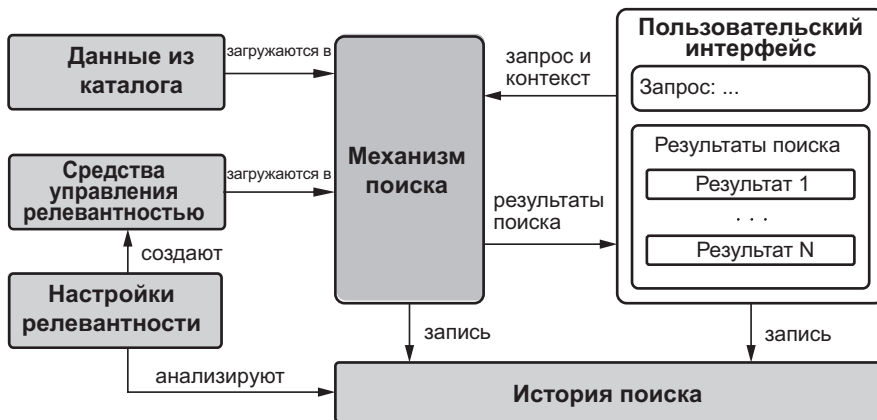


Рис. 4.1. Пример среды поиска для продвижения

- Механизм поиска и прикладной интерфейс фиксируют взаимодействия с пользователями. История взаимодействий обычно включает несколько метрик и параметров, от поисковых запросов до времени прокрутки страницы. Эта информация может использоваться процессом настройки релевантности, который корректирует параметры управления релевантностью для достижения целей.

Среда, которую мы определили выше, ориентирована исключительно на поиск по запросам в свободном текстовом формате, когда пользователь вводит некоторый текст в форму и получает ранжированный список результатов. В действующей поисковой системе, такой как веб-сайт электронной коммерции, эта базовая функциональность часто расширяется множеством дополнительных инструментов и возможностей, включая фильтры, автодополнение запроса и параметры сортировки результатов. Основная функциональность службы поиска почти полностью зависит от понятия релевантности, то есть от разграничения релевантных и нерелевантных элементов. Большая часть этой главы будет посвящена обзору многочисленных аспектов этого вопроса. И для начала мы посмотрим, как можно измерить актуальность и связать ее с экономическими целями.

4.2. Бизнес-цели

Часто можно услышать, что главная цель поисковой службы — понять намерения пользователя и вернуть результаты, соответствующие этим намерениям. В целом это верное утверждение (и трудно реализуемое на практике), однако такой взгляд на цель весьма ограничен и не отражает многих критериев, важных для хорошей поисковой службы. Мы можем попытаться определить всеобъемлющую основу, используя в качестве отправной точки основное уравнение прибыли. Рассмотрим онлайн-магазин, продающий определенный ассортимент товаров, общую прибыль которого можно выразить как

$$G = \sum_j q_j (p_j - v_j), \quad (4.1)$$

где j перебирает все продукты, q представляет проданное количество, p — цену, а v — переменные затраты, которые могут включать оптовую цену и затраты на доставку. Формально мы не можем связать уравнение 4.1 с таким неосознаваемым понятием, как намерение пользователя, но мы можем сделать несколько эвристических предположений, которые помогут установить эту связь и хорошо зарекомендовали себя на практике.

Во-первых, можно предположить, что все проданные объемы q_j примерно пропорциональны органической релевантности и общей эргономике поисковой службы.

То есть учитывая, что число пользователей фиксировано, плохая релевантность дает относительно большую долю пользователей, не сумевших найти товары, соответствующие их покупательским намерениям, а хорошая релевантность дает относительно низкую долю таких пользователей. Следовательно, релевантность можно рассматривать как простой мультипликативный коэффициент прибыли.

Вторая возможность, очевидно вытекающая из уравнения 4.1, заключается в том, что мы можем попытаться перераспределить объемы q_j для продажи более высокодоходных товаров за счет низкодоходных. Можно использовать тот факт, что намерение покупки, как правило, предполагает некоторую гибкость: клиент может изначально стремиться приобрести один товар, но готов заменить его другим, или может просто просматривать доступные предложения, имея неконкретное намерение, будучи готовым принять первый подходящий вариант. Учитывая фиксированное количество пользователей, мы можем продвигать высокодоходные продукты, передвигая их вверх в результатах поиска, чтобы охватить как можно больше потенциальных покупателей, и смещая вниз низкодоходные товары, тем самым отрицательно влияя на их объемы продаж. Важно отметить, что можно учитывать не только прибыль (разницу между стоимостью и продажной ценой), но и затраты на альтернативы, связанные с товаром. Например, торговцы модной одеждой часто практикуют распродажи в конце сезона, чтобы освободить место для новой коллекции, поэтому продвижение устаревающих товаров также может стать целью.

Мы уже говорили, что релевантность напрямую связана с общей прибылью, так как влияет на количество конверсий. Однако мы не должны забывать о негативных последствиях, связанных с релевантностью и дефектами эргономики. С точки зрения пользователя процесс поиска включает несколько шагов, таких как ввод начального запроса, просмотр и проверка предложенных результатов, перестройка запроса и т. д. Плохая релевантность или эргономика могут заставить пользователей повторять поиск много раз, что может негативно сказаться на коэффициенте конверсии, но чаще вызывает у пользователей негативные впечатления и в долгосрочной перспективе уменьшает общее их количество, что в результате дает низкий объем продаж, согласно уравнению 4.1.

Суть в том, что цели службы поиска и ее качество можно рассматривать, по крайней мере, с трех следующих точек зрения:

1. Релевантность.
2. Средства управления продвижением.
3. Эргономика и удовлетворенность клиента.

В следующих разделах мы обсудим каждый из этих трех аспектов, а затем углубимся в детали реализации и посмотрим, как служба поиска может достичь этих целей.

4.2.1. Метрики релевантности

Релевантность результатов поиска можно определить как меру, насколько результаты соответствуют намерениям, с которыми пользователь выполняет поиск. *Цель поиска*, которую в информационном поиске называют также *информационной потребностью*, в большинстве приложений нельзя формализовать полностью. Это особенно верно в сфере продвижения товаров, поэтому стандартным способом измерения релевантности является экспертная оценка результатов поиска и классификации каждого найденного элемента как релевантного или нерелевантного. Например, релевантность результатов поиска по запросу *лечение боли в горле* можно субъективно определить по наличию в них средств лечения этого симптома, начиная от горячего чая до лекарств, а не просто набора элементов, содержащих слова из запроса. Пока будем предполагать, что оценки релевантности задаются экспертами из команды поддержки службы поиска, но, как мы увидим далее в этой главе, поисковая система может собирать и анализировать определенные метрики, характеризующие поведение пользователей, такие как частота щелчков мышью на элементах в результатах, и оценивать релевантность автоматически.

Рассмотрим идеальный случай, когда требуется оценить одну пару запрос/ответ, и общее количество элементов достаточно мало, чтобы их можно было классифицировать вручную. Определим следующие три значения: D — общее количество элементов, релевантных для данного намерения, индексированного в системе, S — количество элементов в результатах поиска и R — количество релевантных элементов. Взаимосвязь между этими тремя значениями показана на рис. 4.2. Качество результатов поиска можно измерить с помощью двух метрик: *точность* (precision) и *полнота* (recall), определяемых следующим образом:

$$\text{точность} = \frac{R}{S}, \quad (4.2)$$

$$\text{полнота} = \frac{R}{D}. \quad (4.3)$$

Обычно для описания результатов или метода поиска требуются обе эти метрики. С одной стороны, полнота отражает полноту результатов поиска, независимо от количества результатов, поэтому всегда можно достичь максимально возможной полноты **1,00**, возвращая всю коллекцию элементов. С другой стороны, точность отражает плотность релевантных элементов в результатах и ничего не говорит о релевантных элементах, которые не были выбраны.

Различие между этими двумя метриками, однако, не означает, что они независимы. Во-первых, еще раз взгляните на рис. 4.2. Судя по рисунку, мы можем изменить полноту от **0** до **1**, растянув прямоугольник результатов поиска по вертикали;

точность при этом останется постоянной. В реальных данных такое поведение практически никогда не наблюдается.

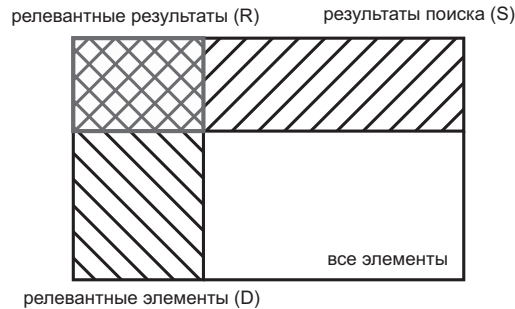


Рис. 4.2. Связь релевантных элементов, результатов поиска и релевантных результатов

Одна из основных причин заключается в том, что мы описываем элементы и определяем запросы с помощью измерений, которые не всегда совпадают с формой набора элементов, определяемой целью поиска. Рассмотрим пример, изображенный на рис. 4.3. У ретейлера имеется большая коллекция обуви, которая описывается с помощью таких свойств, как цена и категория. Пользователь, ищущий *доступную качественную обувь*, может считать в целом релевантными предложения, разбросанные по диагонали и представляющие обувь от дешевых ботинок до дорогих сандалий. Однако поисковая система не сможет повторить эту форму. Буквально интерпретируя критерии *доступный* и *качественный*, она сможет достичь высокой точности, но даст относительно низкую полноту, вернув, например, посредственные кроссовки. Ослабление критериев увеличит полноту, но и добавит в результаты нерелевантные предметы, например дорогие туфли, что ухудшит точность¹.

Этот шаблон почти всегда присутствует в приложениях поиска, поэтому нам часто приходится выбирать между высокоточными методами поиска с низкой полнотой и неточными альтернативами с высокой полнотой. Поиск для продвижения сильно смещен в сторону высокоточных методов, потому что основная его цель состоит в том, чтобы дать пользователю разумное количество релевантных результатов, которые можно быстро просмотреть.

¹ Эта проблема не является чем-то характерным только для поиска и часто возникает в машинном обучении, особенно в приложениях глубокого обучения. Например, набор фотографических изображений представляет собой весьма «замысловатые» области в пространстве всех возможных двумерных матриц. Такие множества, встроенные в многомерные пространства, называются *многообразиями*.

Базовые метрики точности и полноты обеспечивают полезное концептуальное представление релевантности, но как количественные показатели имеют много ограничений. Во-первых, точность и полнота — это метрики для множеств, которые нельзя напрямую применить к ранжированным результатам поиска. Это ограничение имеет решающее значение в поиске для продвижения, который стремится дать пользователю несколько ценных результатов, отсортированных по релевантности. Один из возможных подходов к учету ранжирования — просмотр элементов в результатах сверху вниз, вычисление точности и полноты в каждой точке и построение *кривой точность/полнота*. Этот процесс изображен на рис. 4.4. Допустим, что у нас есть 20 позиций, из которых 5 — релевантны. Набор результатов начинается с релевантного элемента, поэтому точность в этой точке равна 1,00, а полнота 1/5. Следующие два элемента нерелевантны, поэтому полнота остается постоянной, а точность сначала падает до 1/2, а затем до 1/3. Продолжая этот процесс, в двадцатом элементе результатов получаем точность 1/4 и полноту 1,00. То есть кривая точность/полнота — неровная и обычно имеет впадину, выгнутую вниз.



Рис. 4.3. Связь между точностью и полнотой

Кривая точность/полнота дает удобную возможность анализа качества поиска для одного запроса, но часто требуется компактная метрика, выражающая общее качество службы поиска в виде одного числа. Одно из стандартных решений — найти *усредненную среднюю точность* (Mean Average Precision, MAP), то есть сначала найти среднюю точность по каждому релевантному элементу, а затем усреднить это значение по всем запросам, использованным для оценки. Если число запросов равно Q , число релевантных элементов для запроса q равно R_q , а точность k -го релевантного элемента равна P_{qk} , тогда

$$\text{MAP} = \frac{1}{Q} \sum_{q=1}^Q \frac{1}{R_q} \sum_{k=1}^{R_q} P_{qk}. \quad (4.4)$$

Например, на рис. 4.4 показана усредненная средняя точность для одного запроса — это пять точных цифр для каждого из пяти релевантных элементов:

$$\text{MAP} = \frac{1}{5} (1,00 + 0,50 + 0,60 + 0,57 + 0,33) = 0,6. \quad (4.5)$$

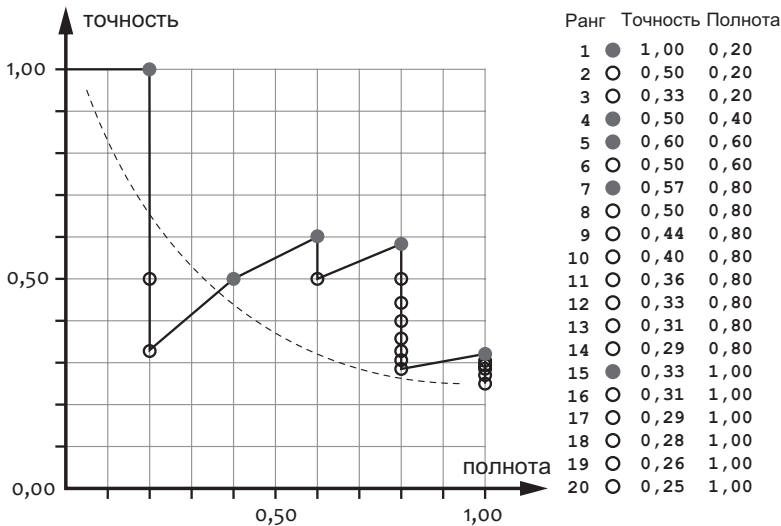


Рис. 4.4. Кривая точность/полнота. Закрашенные кружки соответствуют релевантным элементам в наборе результатов, а незакрашенные — нерелевантным

Обычно в веб-приложениях или приложениях поиска для продвижения нет возможности оценить и перечислить все релевантные результаты для данного запроса, потому что общее их количество огромно; но даже при наличии такой возможности в этом нет никакого смысла, потому что едва ли найдется хоть один пользователь, кому было бы интересно тратить время на чтение всех результатов. Следовательно, приведенная усредненная средняя точность часто вычисляется *не для всех* элементов, относящихся к запросу, а только для релевантных в наборе результатов фиксированного размера [Manning et al., 2008].

Популярной альтернативой усредненной средней точности является метрика *взвешенной накопленной релевантности* (Discounted Cumulative Gain, DCG), которая никак не учитывает понятия точности и полноты и задает центральное место в рейтинге [Järvelin и Kekäläinen, 2000]. Рассмотрим результаты поиска,

содержащего K элементов, каждый из которых градуирован значением релевантности R . То есть у нас имеется k классов релевантности R_k . Значения релевантности могут быть бинарными (единица для релевантных и ноль для нерелевантных элементов), непрерывными или дискретными с несколькими уровнями. *Накопленная релевантность* (Cumulative Gain, CG) результата определяется как сумма оценок:

$$CG = \sum_{k=1}^K R_k. \quad (4.6)$$

Накопленная релевантность напоминает точность, если учесть, что величина K фиксирована, но позволяет различать документы по их полезности, устанавливая многоуровневые оценки R . Именно поэтому данную метрику называют накопленной релевантностью — она пытается оценить полезность результатов поиска. Недостаток накопленной релевантности состоит в том, что эта метрика не учитывает порядок элементов, поэтому изменение порядка не влияет на величину оценки релевантности. Это можно исправить введением штрафа за позиции релевантных результатов, снижающего оценку релевантности пропорционально позиции в списке. Это ведет нас к понятию метрики взвешенной накопленной релевантности (DCG), которая в качестве понижающего веса использует логарифм позиции:

$$DCG = R_1 + \sum_{k=2}^K \frac{R_k}{\log_2(k)}. \quad (4.7)$$

Однако на практике чаще используется немного другое определение DCG, которое особо выделяет релевантные элементы, присваивая им экспоненциально высокие веса [Burges et al., 2005]:

$$DCG = \sum_{k=2}^K \frac{2^{R_k} - 1}{\log_2(k+1)}. \quad (4.8)$$

Величина DCG, рассчитанная по формуле 4.8, зависит от количества результатов K . Для сравнения метрик DCG, полученных для разных запросов, необходимо их нормализовать. Это можно сделать, вычислив максимально возможное значение DCG, называемое идеальной DCG, и разделив фактическое значение DCG на идеальное, получив в результате нормализованное значение DCG (NDCG):

$$NDCG = \frac{DCG}{Ideal\ DCG}. \quad (4.9)$$

Идеальное значение DCG можно оценить, отсортировав результаты по классам релевантности и применив формулу 4.8 для вычисления соответствующего значения DCG. Следовательно, величина NDCG равна единице для идеального ранжирования. Рассмотрим для примера список результатов поиска с шестью элементами,

которые оцениваются экспертом по шкале от 0 до 4, где значение 0 присваивается нерелевантным результатам, значение 4 — наиболее релевантным:

$$4, 3, 4, 2, 0, 1. \quad (4.10)$$

Величина DCG, рассчитанная в соответствии с формулой 4.8, для данного результата равна 28.56. Идеальный порядок результатов в этом случае

$$4, 4, 3, 2, 1, 0, \quad (4.11)$$

и соответствующее значение идеальной величины DCG равно 29,64. Следовательно, оценка NDCG для результатов в списке 4.10 будет равна $28,56/29,64 = 0,96$.

4.2.2. Средства управления продвижением

Средства управления продвижением — это инструменты, позволяющие мерчандайзерам и другим бизнес-пользователям влиять на результаты поиска в соответствии с потребностями бизнеса, не охваченными органической релевантностью. Однако граница между параметрами управления продвижением и релевантностью очень размыта, поскольку многие бизнес-правила можно рассматривать как улучшения релевантности на основе содержимого, и, наоборот, многие методы управления релевантностью можно рассматривать как бизнес-правила, улучшающие результаты поиска внедрением некоторых знаний о предметной области. Например, мерчандайзер может создать триггер, направляющий всех пользователей, вводящих запрос *утепленные куртки*, в категорию, курируемую вручную и созданную специально для сезонной распродажи курток. С одной стороны, это действие направлено на достижение бизнес-цели по продвижению товаров. С другой — можно утверждать, что такая категория, курируемая вручную, лучше соответствует целям, которые преследует пользователь, чем результаты стандартного поиска. Большинство поисковых систем предоставляют богатый набор средств управления продвижением, который может включать следующие возможности:

ПОВЫШЕНИЕ И ПОНИЖЕНИЕ. Как мы увидим в следующих разделах, органическая релевантность обычно вычисляется путем сопоставления различных свойств запросов и элементов, смешивания полученных оценок и ранжирования элементов по окончательной оценке. Мерчандайзер может скорректировать или переопределить эту логику, изменяя оценки релевантности для продвижения желательных товаров и понижения продаж нежелательных. В мире маркетинга этот прием часто называют *повышением и понижением*. Управление повышением и понижением часто можно выразить в виде формулы оценки, которая смешивает различные свойства элемента. Этот подход можно проиллюстрировать на примере

повышения новых (*newness*), уцененных (*discount*) или высокорейтинговых (*rating*) товаров и понижения товаров, которые не имеют этих свойств:

$$\text{score} = 0,2 \times \text{новые} + 0,4 \times \text{уцененные} + 0,4 \times \text{высокорейтинговые}, \quad (4.12)$$

с учетом того, что каждому товару приписывается соответствующая новизна (*newness*), величина скидки (*discount*) и рейтинг (*rating*), измеренные по некоторой шкале. Оценку, рассчитанную по такой формуле, можно использовать взамен оценки релевантности, или суммировать, или умножать две оценки.

ФИЛЬТРАЦИЯ. Основной целью фильтрации является устранение нежелательных элементов из результатов поиска. В качестве примера фильтрации можно привести удаление товаров, отсутствующих на складе, и удаление нерелевантных товаров, появившихся в результатах поиска из-за проблем с данными или оценкой.

ФИКСИРОВАННЫЕ РЕЗУЛЬТАТЫ. Порой трудно добиться определенного порядка элементов, используя формулу повышения и понижения, поэтому мерчандайзер должен быть готов поместить отобранный вручную набор элементов в верхнюю часть результатов. Внедрение таких фиксированных элементов часто инициируется определенными ключевыми словами в запросе.

ПЕРЕНАПРАВЛЕНИЕ. Прием перенаправления сродни приему внедрения фиксированных результатов, но в отличие от последнего полностью заменяет органические результаты поиска, перенаправляя пользователя в категорию товаров или содержанию, курируемому вручную, например, в интерактивный журнал моды.

ГРУППИРОВКА ТОВАРОВ. Эффективное использование экранного пространства является важной целью поиска для продвижения товаров. Важно не только вернуть пользователям релевантные результаты, но и представить имеющийся ассортимент наилучшим образом, учитывая, что экранное пространство ограничено. Например, иногда полезно заменить тесно связанные товары или их варианты, такие как различные размеры и цвета одной и той же модели джинсов, одним представителем категории, чтобы освободить больше места для других моделей и избежать загромождения результатов поиска аналогичными элементами.

С экономической точки зрения, некоторые средства управления продвижением можно рассматривать как методы сегментации рынка. Рассмотрим пример продвижения высокодоходных товаров или предметов роскоши: по сути это попытка сегментировать клиентов по их чувствительности к ценам, то есть клиенты, нечувствительные к цене, тратят меньше усилий, выбирая высокодоходные товары в верхней части результатов, а клиентам, чувствительным к ценам, приходится тратить больше времени и более скрупулезно просматривать страницы с результатами, чтобы найти лучшую цену.

4.2.3. Метрики качества службы поиска

Релевантность результатов поиска, оцененная экспертом и измеренная с помощью таких метрик, как NDCG, не гарантирует приемлемого качества службы поиска. Нам нужно определить метрики, которые можно измерять и контролировать в реальных приложениях, чтобы достичь максимально положительных впечатлений у пользователя и эффективности бизнеса. Качество службы поиска можно связать с алгоритмами релевантности, качеством данных, эргономичностью пользовательского интерфейса и надежностью технической реализации. Рассмотрим несколько ключевых показателей качества, которые часто используются для оценки служб поиска:

КОЭФФИЦИЕНТ КОНВЕРСИИ. На сегодняшний день коэффициент конверсии является наиболее важным показателем эффективности поиска для продвижения. Его можно определить как отношение числа сеансов пользователей, использовавших службу поиска и сделавших покупку, к общему числу сеансов использования службы поиска. В этом смысле сеанс пользователя обычно эквивалентен веб-сеансу. Коэффициент конверсии является крайне важным показателем, поскольку напрямую связан с доходами и возможностью найти нужные товары.

ПРОЦЕНТ ПЕРЕХОДОВ. Отношение числа пользователей, щелкнувших на конкретном результате, к общему числу пользователей, воспользовавшихся поиском, является важным показателем релевантности результатов.

ВРЕМЯ ПРЕБЫВАНИЯ НА СТРАНИЦЕ С ОПИСАНИЕМ ТОВАРА. Высокий процент переходов по ссылке в целом позитивный показатель, но большое количество пользователей, которые, зайдя на страницу с описанием, быстро возвращаются обратно, может указывать на плохую релевантность или проблемы с эргономикой, из-за которых пользователь не может распознать товары по их кратким описаниям в списке результатов.

ЧАСТОТА ИЗМЕНЕНИЯ ЗАПРОСА. Если пользователь изменяет запрос несколько раз, высока вероятность того, что он не может получить удовлетворительные результаты.

ЧАСТОТА ЛИСТАНИЯ СТРАНИЦ. Высокая частота перемещения между страницами и щелчков на результатах с низким рейтингом может указывать на проблемы с релевантностью.

КОЭФФИЦИЕНТ УДЕРЖАНИЯ. Доля пользователей, продолжающих регулярно пользоваться поиском. Коэффициент удержания обычно рассчитывается для определенного периода времени, например недели или месяца, как

$$\text{коэффициент удержания} = \frac{E - N}{S}, \quad (4.13)$$

где E — число постоянных пользователей в конце периода, N — число новых пользователей, привлеченных в этот же период, и S — число пользователей в начале периода.

ЗАДЕРЖКА РЕЗУЛЬТАТОВ. Время, необходимое для обработки поискового запроса и возврата результатов, оказывает большое влияние на впечатления пользователя. Многие ретейлеры и компании, предлагающие услуги веб-поиска, дают впечатляющую статистику по этому вопросу. Например, Amazon сообщила, что увеличение времени загрузки страницы на каждые 100 миллисекунд приводит к уменьшению продаж на 1 %, а Walmart отметила, что уменьшение задержки на каждую секунду увеличивает конверсию на 2 % [Kohavi and Longbotham, 2007; Crocker et al., 2012].

Эти показатели можно разбить по измерениям, таким как маркетинговые каналы (например, мобильный, настольный или планшетный), для точной настройки релевантности. Мы еще вернемся к вопросу настройки релевантности, включая ручную и автоматическую настройку, в разделе 4.7, но сначала познакомимся поближе с особенностями оценки релевантности.

4.3. Строительные блоки: соответствие и ранжирование

Проблему релевантности поиска можно рассматривать как проблему классификации, потому что она нацелена на различение релевантных и нерелевантных элементов. В то же время это очень специфический случай классификации из-за ориентации на текстовые данные и ранжирование. Эти особенности позволяют использовать очень эффективные эвристические методы, способные достичь высочайшей релевантности без обучения классификатора с помощью методов машинного обучения. Подход на основе машинного обучения также возможен (и мы обсудим его в следующих разделах), но для большинства приложений поиска для продвижения достаточно базовых методов сопоставления и ранжирования, которые, к тому же, образуют надежную методологию для проектирования признаков в случаях, когда применяются методы машинного обучения. В этом разделе мы рассмотрим базовые методы поиска, которые потом можно использовать для создания более сложных и комплексных решений релевантности, а также для настройки параметров в машинном обучении.

В общих чертах поиск можно описать как вычисление некоторой метрики сходства между элементом и запросом с последующим ранжированием элементов в соот-

ветствии с этой метрикой и исключением нерелевантных элементов из результатов. Как и в других задачах классификации, для этого требуется получить представление элемента, запроса и, при необходимости, другой контекстной информации, такой как сведения из профиля пользователя, в виде объектов, а затем вычислить одну или нескольких оценок, которые далее мы будем называть *сигналами*, указывающих, насколько хорошо объект элемента соответствует объекту запроса. После этого сигналы объединяются, и принимается окончательное решение о включении элемента в список результатов (*соответствие*) и его позиции в списке (*ранжирование*). Рисунок 4.5 иллюстрирует этот поток.

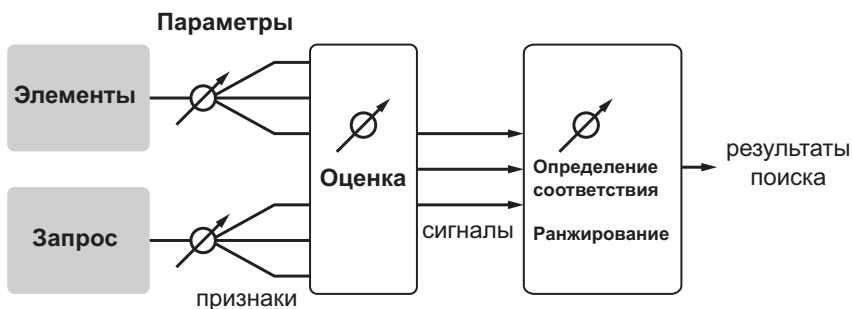


Рис. 4.5. Обобщенный поток поиска и его основные задачи

Как следствие, на каждом из этих этапов разработчик службы поиска должен использовать множество параметров, определяющих, как извлекать признаки из исходных данных, как сопоставлять признаки для получения сигналов и как смешивать сигналы для достижения хорошей релевантности. Для начала рассмотрим несколько основных методов, а затем постепенно будем увеличивать сложность, объединяя блоки и добавляя в уравнение новые переменные.

4.3.1. Лексическое соответствие

В поиске для продвижения элементы обычно представлены довольно сложными сущностями с несколькими текстовыми и числовыми атрибутами, такими как имя, описание, цена и бренд. Данные также могут содержать значительный объем структурной информации, такой как иерархия категорий и различные варианты размер/цвет одного логического продукта. Давайте оставим пока все эти сложности в стороне и рассмотрим простой случай ретейлера, который представляет каждый продукт как документ с одним полем описания, содержащим простой текст, например:

Product 1

Description: Pleated black dress. Lightweight look
for the office.¹

Product 2

Description: Fiery red dress. A black ribbon
at the waist.²

Самое простое, что можно сделать для поиска по таким документам, — разбить описания на слова и разрешить поисковые запросы только по одному слову, то есть продукт будет включен в результаты, только если его описание содержит слово из запроса. Процесс разбиения текста на слова или другие элементы, такие как фразы, называется *лексемизацией*, а выходные данные — в нашем случае слова — называются лексемами. В английском языке (как и в русском. — *Примеч. пер.*) лексемизация обычно выполняется с использованием пробелов и знаков препинания в качестве разделителей, поэтому для документов выше будут созданы следующие списки лексем:

Product 1: [Pleated], [black], [dress], [Lightweight],
[look], [for], [the], [office]

Product 2: [Fiery], [red], [dress], [A], [black]
[ribbon], [at], [the], [waist] (4.14)

Следовательно, запросу *black* будут соответствовать оба продукта, а запросу *red* — только второй. Очевидно, что этот метод дает только возможность определения соответствия и не ранжирует элементы в результатах.

Несмотря на простоту и ограниченность, метод сопоставления лексем иллюстрирует основные принципы управления поиском, рассмотренные выше. Во-первых, каждую лексему можно рассматривать как отдельный признак, который может присутствовать или отсутствовать в описании продукта. Процесс лексемизации служит примером проектирования признаков. Слова действительно являются достаточно хорошими признаками, потому что несут сильный сигнал о типе продукта, таком как обувь, и его свойствах, таких как черный (*black*) цвет. Во-вторых, сопоставление лексем — это способ получения сигналов о корреляции между признаками продукта и запроса. Наконец, сигналы от всех лексем объединяются вместе для принятия окончательного решения о соответствии или несоответствии. Этот поток представлен на рис. 4.6.

¹ Плиссированное черное платье. Хорошо подходит для работы в офисе.

² Огненно-красное платье. С черным поясом на талии.

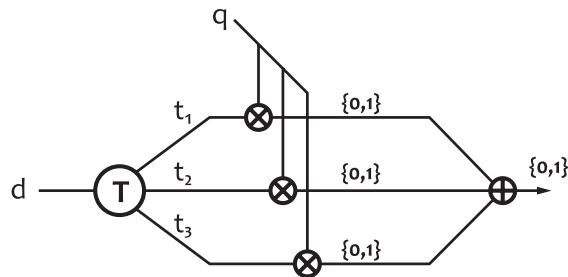


Рис. 4.6. Иллюстрация потока сопоставления лексем, включающего лексемизацию, оценку и смешивание сигналов. Описание продукта и его лексем обозначены как d и t_i соответственно. Запрос обозначен как q

4.3.2. Логический поиск и поиск по фразе

Основным недостатком сопоставления простых лексем является невозможность обработки более осмысленных запросов, включающих более одного критерия. Первое улучшение, которое поможет обойти это ограничение, — включение поддержки *логических запросов*, которая позволит связать несколько лексем с помощью логических операций, а именно AND (И), OR (ИЛИ) и NOT (НЕ). Например, следующему логическому запросу

`dress AND red`

будет соответствовать только второй продукт из примера 4.14, тогда как запросу

`dress AND (red OR black)`

будут соответствовать оба продукта. Логические запросы не учитывают позиции лексем в тексте, и, соответственно, их можно рассматривать как объединение нескольких запросов сопоставления лексем.

Второй важной особенностью, расширяющей простое сопоставление лексем, является поддержка *фразовых запросов*. Фразовый запрос — это запрос, выполняющий поиск документов с последовательностью лексем, которые следуют друг за другом, в отличие от логического запроса, который ищет документы, содержащие отдельные лексем, независимо от их порядка и местоположения в тексте. Для обозначения запросов с фразами и подзапросов мы будем использовать квадратные скобки. Например, следующему запросу будет соответствовать первый продукт из примера 4.14, но не второй:

`[black dress]`

Этот результат имеет более высокую точность и низкую полноту, чем результат логического запроса *black AND dress*, которому соответствуют оба продукта. Логические и фразовые запросы вместе являются очень мощными инструментами контроля релевантности и управления компромиссом между точностью и полнотой. Язык запросов, непосредственно поддерживающий логические выражения, часто является хорошим решением для экспертного поиска, когда пользователи готовы изучать и использовать расширенные функции поиска, но его использование в поиске для продвижения ограничено из-за неготовности рядового пользователя. В следующих разделах мы обсудим, как текстовый поиск может использовать преимущества сложных логических и фразовых запросов.

4.3.3. Нормализация и стемминг

Нетрудно заметить, что разделение текста на лексемы дает не самые оптимальные результаты с точки зрения соответствия. В естественном языке слова могут иметь разные формы написания, которые почти неразличимы для целей поиска. Некоторые слова вообще не несут никакой значимой информации и генерируют помеху. То есть нам необходимо предусмотреть нормализацию исходных лексем для создания более чистого словаря. Такие нормализованные лексемы обычно называются *термами*, или *терминами*.

Нормализация — сложный процесс, который обычно включает несколько шагов для учета различных свойств и явлений естественного языка. Рассмотрим пример, иллюстрирующий ключевые преобразования, начиная со следующего исходного описания продукта:

Maison Kitsuné Men's Slim Jeans. These premium jeans come in a slim fit for a fashionable look.¹

Первым шагом является нормализация набора символов, потому что поисковый запрос может вводиться с диакритическими знаками или без них, и это различие обычно не означает разных целей поиска. Лексемизация текста и его приведение к стандартному набору символов дают нам следующие лексемы:

```
[Maison] [Kitsune] [Men's] [Slim] [Jeans] [These]  
[premium] [jeans] [come] [in] [a] [slim] [fit]  
[for] [a] [fashionable] [look]
```

Вторая проблема, с которой мы сталкиваемся, — наличие символов нижнего и верхнего регистров, которые также неразличимы в большинстве случаев. Обычно для

¹ Maison Kitsuné (название бренда) — мужские узкие джинсы. Модные облегающие джинсы премиум-класса. — *Примеч. пер.*

решения этой проблемы все символы приводятся к нижнему регистру, в результате чего мы получаем следующий результат для нашего примера:

```
[maison] [kitsune] [men's] [slim] [jeans] [these]  
[premium] [jeans] [come] [in] [a] [slim] [fit]  
[for] [a] [fashionable] [look]
```

Третий возможный шаг — исключение часто встречающихся лексем, таких как *and*, *to*, *the* и *will*, потому что они присутствуют в большинстве текстов и не несут никакой определенной информации об элементе. Такие лексемы обычно называются *стоп-словами*. Выполнив такое исключение в данном примере, получим

```
[maison] [kitsune] [men's] [slim] [jeans] [premium]  
[jeans] [come] [slim] [fit] [fashionable] [look]
```

Исключение стоп-слов может иметь как положительные, так и отрицательные последствия. С одной стороны, этот шаг может положительно повлиять на некоторые методы сопоставления и ранжирования, которые мы обсудим позже, потому что часто встречающиеся лексемы, не несущие смысловой нагрузки, могут исказить некоторые метрики, используемые в расчетах релевантности. С другой стороны, удаление стоп-слов может привести к потере существенной информации и препятствовать поиску по определенным фразам. Например, удаление стоп-слов не позволит выполнить поиск по фразе *to be, or not to be* (быть или не быть) или отличить *new* (новое) от *not new* (не новое). Стоп-слова также могут разрушать семантические связи между сущностями, делая неразличимыми объекты *on the table* (на столе) и *under the table* (под столом).

Четвертый стандартный метод нормализации — *стемминг*. В большинстве естественных языков слова могут менять форму в зависимости от числа (*платье* и *платья*), времени (*вижу* и *видел*), владения (*мужчина* и *мужской*) и других факторов. Стемминг — это процесс выделения основы слова с целью ликвидации различий, не влияющих на цели поиска. Проблема стемминга сопряжена со сложностями из-за многочисленных исключений и особых случаев в естественных языках. Существует множество методов стемминга, основанных на правилах или словарях, обладающих своими достоинствами и недостатками. Одно из популярных семейств методов основано на так называемом стеммере Портера [Porter, 1980]. Он представляет несколько групп правил преобразования суффиксов и условий, помогающих убедиться, что слово достаточно длинное, чтобы его можно было сократить (см. табл. 4.1). Применяв стемминг к нашему примеру описания продукта, мы получим окончательный набор термов, более компактный и сфокусированный на ключевых особенностях, чем исходный текст:

```
[maison] [kitsun] [men] [slim] [jean] [premium]  
[jean] [com] [slim] [fit] [fashion] [look]
```

Таблица 4.1. Пример правил, используемых в стеммере Портера. Все правила в этом примере требуют, по крайней мере, одного перехода от гласной к согласной перед суффиксом, поэтому второе правило применяется к слову *conditional*, но не применяется к *rational*

Правило	Пример
...ational → ate	relational → relate
...tional → tion	conditional → condition rational → rational
...ful → <i>ничего</i>	hopeful → hope
...ness → <i>ничего</i>	goodness → good
...izer → ize	digitizer → digitize

Один и тот же набор алгоритмов нормализации обычно применяется и к запросу, и к документам, чтобы отобразить все лексемы в одно и то же пространство термов. Например, запросу *Fashionable* (модно) будет соответствовать продукт, содержащий в описании слово *fashioned* (модный), потому что оба слова отобразятся в терм *fashion* (мода).

4.3.4. Ранжирование и модель векторного пространства

Сопоставление с помощью логических и фразовых запросов позволяет найти набор элементов, удовлетворяющих критериям поиска. Однако количество соответствующих и потенциально релевантных элементов часто превышает относительно небольшое количество результатов, которое готов просмотреть средний пользователь, поэтому порядок, в каком элементы будут представлены пользователю, становится критически важным. Нам нужно определить строительный блок, который будет ранжировать элементы в соответствии с их релевантностью.

Ранжирование не улучшает глобальных свойств точности/полноты базового сопоставления, но его можно считать своеобразным трюком, улучшающим точность/полноту в смысле локальных или воспринимаемых качеств. С одной стороны, ранжирование повышает точность верхних результатов, повышая релевантность элементов, но в то же время не удаляет элементы из результатов, что обеспечивает сохранение полноты.

Первым шагом на пути к ранжированию внимательнее рассмотрим логические запросы и определим их оценочный потенциал. Прежде всего отметьте, что документы и запросы можно представить в виде бинарных векторов, в которых каждый элемент указывает, содержит ли документ или запрос определенный терм. Другими

словами, элемент, соответствующий определенному терму, равен единице, если документ или запрос содержит этот терм, и нулю в противном случае. Если общее число различных термов во всех документах коллекции равно n , то каждый документ или запрос является бинарным вектором с n элементами. Легко увидеть, что логический запрос можно выразить как скалярное произведение между вектором запроса \mathbf{q} и вектором документа \mathbf{d} . Напомню, что скалярное произведение двух векторов задается как

$$\mathbf{d} \cdot \mathbf{q} = \sum_{i=1}^n q_i d_i, \quad (4.15)$$

а евклидова норма определяется как

$$\|\mathbf{d}\| = \sqrt{d_1^2 + \dots + d_n^2}. \quad (4.16)$$

Следовательно, можно сказать, что логический запрос с несколькими термами, связанными с помощью оператора *AND*, эквивалентен следующему условию:

$$\mathbf{d} \cdot \mathbf{q} \geq \|\mathbf{q}\|^2, \quad (4.17)$$

поскольку для всех термов в запросе должны иметься соответствующие термы в документе, скалярное произведение должно быть равно числу элементов в векторе запроса. Логический запрос, в котором термы связаны условием *OR*, эквивалентен условию

$$\mathbf{d} \cdot \mathbf{q} \geq 1, \quad (4.18)$$

поскольку должно иметься хотя бы одно совпадение. Такая интерпретация логических запросов демонстрирует вид внутренней оценки, которая преобразуется в решение о соответствии с применением порога. Уравнения 4.17 и 4.18 также предполагают, что отношение между скалярным произведением и нормой запроса можно использовать как непрерывную меру сходства между документом и запросом. Мы можем пойти еще дальше и спросить, почему не учитывается норма документа? Можно утверждать, что короткий документ, соответствующий условиям запроса, релевантнее длинного документа, соответствующему тому же числу термов. Это можно доказать с вероятностной точки зрения следующим образом. Представьте человека, произносящего речь по какой-то теме: если он использует много релевантных слов в первую минуту, это может служить показателем, что речь действительно посвящена соответствующей теме. В то же время часовая речь, в которой используются те же слова, может быть посвящена широкому разнообразию тем. Метрика, нормализующая скалярное произведение по нормам обоих векторов, называется *косинусным сходством* (cosine similarity):

$$\cos(q, d) = \frac{q \cdot d}{\|q\| \cdot \|d\|}. \quad (4.19)$$

Косинусное сходство, то есть косинус угла между векторами, является удобной метрикой, изменяющейся в диапазоне от нуля до единицы для положительно определенных векторов. Косинусное сходство, равное нулю, означает, что вектор документа ортогонален вектору запроса в пространстве термов, а сходство, равное единице, означает точное соответствие логическому запросу. В отличие от логического запроса, косинусное сходство не требует указывать операции в запросе — запрос и документ рассматриваются как неупорядоченные коллекции термов. Проиллюстрируем эту *модель векторного пространства* на примере.

ПРИМЕР 4.1

Рассмотрим два элемента, имеющие следующие описания (для упрощения предположим, что описания лексемизированы и нормализованы):

Product 1: dark blue jeans blue denim fabric

Product 2: skinny jeans in bright blue

Эти два описания и запрос *dark jeans* представлены в табл. 4.2 в виде бинарных векторов.

Таблица 4.2. Пример представления двух документов и одного запроса в виде бинарных векторов

	dark	blue	jeans	denim	fabric	skinny	in	bright	·
d1	1	1	1	1	1	0	0	0	$\sqrt{5}$
d2	0	1	1	0	0	1	1	1	$\sqrt{5}$
q	1	0	1	0	0	0	0	0	$\sqrt{2}$

Значения сходства между запросом и каждым из документов:

$$\begin{aligned} \cos(q, d_1) &= \frac{1+1}{\sqrt{2}\sqrt{5}} = 0,632, \\ \cos(q, d_2) &= \frac{1}{\sqrt{2}\sqrt{5}} = 0,316. \end{aligned} \quad (4.20)$$

На рис. 4.7 показана связь между документами и запросами в векторном пространстве. Обратите внимание, что косинусное сходство имеет эф-

фективное решение, так как требуется учесть только ненулевые элементы вектора запроса.

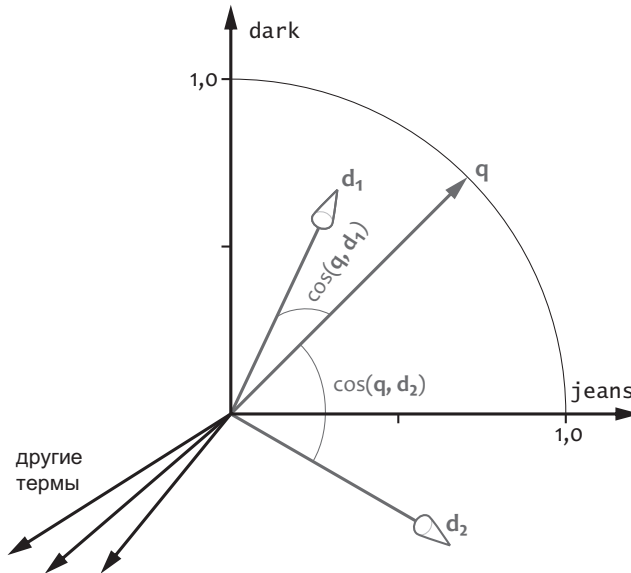


Рис. 4.7. Пример модели векторного пространства и косинусное сходство двух документов и запроса. Векторы документов и запроса нормализованы

4.3.5. Модель оценки $TF \times IDF$

Модель векторного пространства с бинарными векторами имеет два важных недостатка, которые отрицательно влияют на релевантность результатов, ранжированных с помощью этого метода. Во-первых, он не учитывает частоту термов в документе. Можно ожидать, что документы с несколькими вхождениями термов из запроса более релевантны, чем документы, в которых те же термы встречаются только один раз. Во-вторых, некоторые термы могут быть более важными: редко используемые слова часто являются более важными и информативными, чем часто используемые. Например, в описаниях товаров продавец одежды может часто использовать слово *одежда*, поэтому соответствие этому терму не является сильным сигналом релевантности. Стоп-слова, которые мы рассматривали выше, являются крайним случаем этой проблемы.

Первую проблему можно устранить заменой нулей и единиц в векторах документов соответствующими частотами термов в документе. Этот вариант модели

векторного пространства часто называют *моделью мешка слов*. *Частоту термина* (Term Frequency, TF) можно определить как число его вхождений t в документ d , который мы обозначим как $n(t, d)$ — некоторую нелинейную функцию от этого числа. Одно из популярных решений заключается в использовании квадратного корня из числа вхождений:

$$\text{tf}(t, d) = \sqrt{n(t, d)}. \quad (4.21)$$

Например, терм, встречающийся в документе девять раз, будет иметь частоту термина, равную трем. Функция квадратного корня используется для сглаживания оценок документов, имеющих очень большое количество вхождений термов.

Вторую проблему можно решить, вычисляя частоты термов по всей коллекции документов, и с их помощью отличать редкие слова от частых. Один из возможных способов оценки редкости слов — подсчет вхождений термина во всех документах аналогично методу определения частоты термов, но для всей коллекции. Однако этот подход, как известно, дает неоптимальные результаты, потому что несколько документов с многочисленными вхождениями редкого термина могут исказить результаты. Более распространен получил метод подсчета количества документов, содержащих хотя бы одно вхождение данного термина. Эта метрика называется *частотой документа* (document frequency) термина. Соответственно, в качестве показателя редкости термов можно использовать *обратную частоту документа* (Inverse Document Frequency, IDF). Стандартная формула обратной частоты документа (IDF) для термина t выглядит следующим образом:

$$\text{idf}(t) = 1 + \ln \frac{N}{\text{df}(t) + 1}, \quad (4.22)$$

где N — общее число документов в коллекции, а $\text{df}(t)$ — частота документа для термина. По аналогии с частотой термина, чтобы сгладить величину коэффициента для редких членов, используется логарифмическая функция.

Частота термина и обратная частота документа часто объединяются вместе, то есть элементы вектора документа вычисляются как произведение значений, определяемых уравнениями 4.21 и 4.22:

$$d(i) = \text{tf}(t_i, d) \times \text{idf}(t_i). \quad (4.23)$$

Этот подход, известный как модель $\text{TF} \times \text{IDF}$, получил широкое распространение. Подставляя выражение 4.23 в определения скалярного произведения и евклидовой нормы, получаем следующие формулы, которые можно использовать для вычисления оценки косинусного сходства запроса q и документа d в рамках модели $\text{TF} \times \text{IDF}$:

$$\mathbf{q} \cdot \mathbf{d} = \sum_{t \in q} \text{tf}(t, d) \cdot \text{idf}(t) \times \text{tf}(t, q) \cdot \text{idf}(t), \quad (4.24)$$

$$\|\mathbf{q}\| = \sqrt{\sum_{t \in q} [\text{tf}(t, q) \cdot \text{idf}(t)]^2}, \quad (4.25)$$

$$\|\mathbf{d}\| = \sqrt{\sum_{t \in d} [\text{tf}(t, d) \cdot \text{idf}(t)]^2}. \quad (4.26)$$

Косинусное сходство, вычисленное с помощью формул 4.24–4.26, является одним из наиболее широко используемых методов оценки. Эти формулы, однако, не являются надежным стандартом и во многих реализациях поисковых систем можно найти множество различных вариантов. Примером могут служить три следующие поправки, хорошо зарекомендовавшие себя на практике:

1. Норма документа, определяемая уравнением 4.26, нормализует векторы всех документов, приводя их к единичной длине. Однако во многих практических приложениях более короткие документы часто оказываются более релевантными, если они содержат одинаковое количество совпадений термов и равную частоту термов. Это обстоятельство можно учесть, заменив стандартную норму документа нормой, пропорциональной общему количеству термов $n(d)$ в документе:

$$L_d(d) = \sqrt{\sum_{t \in d} 1} = \sqrt{n(d)}. \quad (4.27)$$

2. Все термы в запросе могут считаться одинаково значимыми и обрабатываться независимо, даже если они повторяются. Следовательно, частота терма $\text{tf}(t, q)$ всегда равна единице. Это позволяет переопределить норму запроса следующим образом:

$$L_q(q) = \sqrt{\sum_{t \in q} \text{idf}(t)^2}. \quad (4.28)$$

3. Оценка документа в модели $\text{TF} \times \text{IDF}$ зависит от количества слов, соответствующих запросу, поскольку пропущенные слова обнуляют соответствующие термы в скалярном произведении. Можно утверждать, что штраф за пропущенные слова должен быть еще больше, с этой целью можно ввести дополнительный коэффициент, называемый координационным фактором. *Координационный фактор* $c(q, d)$ определяется как отношение числа часто встречающихся термов в запросе и документе к общему числу термов в запросе. Например, запрос *black skinny jeans* и документ *black jeans* — координационный фактор, равный двум третьим.

Подставив все эти поправки в определение косинусного сходства, получим следующую итоговую формулу оценки:

$$\text{score}(q, d) = \frac{c(q, d)}{L_d(d) \cdot L_q(q)} \sum_{t \in q} \text{tf}(t, d) \cdot \text{idf}(t)^2. \quad (4.29)$$

Модель TF×IDF является фундаментальным строительным блоком, который мы будем использовать позже для создания более сложных решений оценки. Стоит также отметить, что этот метод разрабатывался как обобщенное решение для поиска относительно длинных текстов, таких как журнальные статьи, и его применение в поиске для продвижения, имеющем дело со структурированными данными, может приводить к непредсказуемым последствиям, из-за чего иногда желательно задействовать альтернативные методы ранжирования, как будет показано позже.

ПРИМЕР 4.2

Завершим обзор методов ранжирования примером вычисления TF×IDF. Используем следующие описания продуктов:

d_1 : dark blue jeans blue denim fabric

d_2 : skinny jeans in bright blue

Применяя формулы 4.21 и 4.22, получаем значения TF и IDF, представленные в табл. 4.3. Теперь оценим соответствие этих продуктов запросу *skinny jeans*. Затем рассчитаем нормы запроса и документов в соответствии с формулами 4.27 и 4.28, используя только что полученные значения TF и IDF.

$$L_d(d_1) = \sqrt{6} = 2,449, \quad L_d(d_2) = \sqrt{5} = 2,236, \quad (4.30)$$

$$L_q(q) = \sqrt{\text{idf}(\text{jeans})^2 + \text{idf}(\text{skinny})^2} = 1,163. \quad (4.31)$$

Таблица 4.3. Пример расчета TF и IDF для двух документов. Последние две строки соответствуют векторным представлениям TF×IDF документов

	dark	blue	jeans	denim	fabric	skinny	in	bright
idf(·)	1,00	0,59	0,59	1,00	1,00	1,00	1,00	1,00
tf(·, 1)	1,00	1,41	1,00	1,00	1,00	0,00	0,00	0,00
tf(·, 2)	0,00	1,00	1,00	0,00	0,00	1,00	1,00	1,00
d1	1,00	0,83	0,59	1,00	1,00	0,00	0,00	0,00
d2	0,00	0,59	0,59	0,00	0,00	1,00	1,00	1,00

Координационный фактор равен 0,50 для первого продукта и 1,00 для второго. Подставляя все нормы и значения $TF \times IDF$ в формулу 4.29, получаем оценку 0,062 для первого продукта и гораздо более высокую оценку 0,520 для второго, что согласуется с интуитивным представлением о большей релевантности второго продукта.

$TF \times IDF$ также зависит от нормализации текста и стемминга. Например, не трудно заметить, что если не выполнить стемминг, следующие документы получат равные оценки $TF \times IDF$ для запроса *dark*:

d_1 : dark darker darkness
 d_2 : dark darker lightness
 d_3 : dark light lightness

Однако в этом контексте первый документ выглядит более релевантным. Стемминг отобразит слова *dark*, *darker* и *darkness* в одну и ту же основу *dark*, благодаря чему первый и второй документы получают более высокие оценки из-за более высокой частоты термов. Кроме того, пользователь, ищущий *darkish shoes* (темные ботинки), не получит никаких результатов без стемминга, что вряд ли будет способствовать появлению у него положительного восприятия.

4.3.6. Оценка с использованием n -грамм

Теперь мы знаем, что модель векторного пространства связана с логическими запросами, и оценку $TF \times IDF$ можно рассматривать как смягченную разновидность логического запроса, которая заполняет пробел между логическими OR- и AND-запросами. Мы можем продолжить рассуждать в том же направлении и определить смягченную версию фразового запроса. С точки зрения ранжирования стандартные логические фразовые запросы являются слишком строгими и требуют наличия совпадений со всеми термами в запросе. Один из возможных способов ослабить это требование — ограничиться совпадениями с *фрагментами* фразы, то есть последовательностями из нескольких термов. Такие последовательности также называют *n -граммами* и могут состоять из двух (биграммы), трех (триграммы) или более термов. Выделение фрагментов можно рассматривать как метод лексемизации, применяемый к документу и к запросу, и полученные n -граммы включать в логический запрос или механизм вычисления оценки $TF \times IDF$. В следующем примере показано, как описание двух продуктов можно разделить на биграммы:

```
black cotton polo shirt: [black cotton]
                        [cotton polo]
                        [polo shirt]
```

```
short sleeve black shirt: [short sleeve]
                           [sleeve black]
                           [black shirt]
```

Механизм оценки $TF \times IDF$ обрабатывает n -граммы так же, как однословные термины, и вычисляет косинусное сходство в векторном пространстве, только при этом каждый элемент вектора соответствует фрагменту фразы, и метрики $TF \times IDF$ вычисляются для фрагментов. Как результат, если применить разбиение на отдельные слова (униграммы), описания получают равные оценки $TF \times IDF$ для запроса *black shirt*, но, если применить разбиение на биграммы, второй продукт получит более высокую оценку, потому что содержит точный фрагмент *black shirt*. Можно утверждать, что оценка с применением биграмм лучше фиксирует семантические связи в тексте: близость слов *polo* и *shirt* в описании первого продукта подчеркивает, что *polo shirt* (рубашка поло) — это главное его свойство, а цвет *black* (черный) лишь уточняющее, тогда как близость *black* и *shirt* в описании второго продукта указывает, что черный цвет (*black*) — это ключевое свойство. Такая способность выделять семантические связи особенно важна для различения составных термов, таких *tuxedo coat* (смокинг) и *sports coat* (спортивная куртка). Использование фрагментов — мощный метод повышения точности поиска, который часто сочетается со стандартной оценкой одиночных слов, о чем мы поговорим чуть ниже.

4.4. Смешивание сигналов релевантности

До сих пор мы рассматривали поиск элементов в простых текстовых описаниях. В поиске для продвижения, как и в большинстве других поисковых приложений, такой простой формат данных встречается редко. Нам почти всегда придется иметь дело со структурированными исходными данными, характеризующими каждый элемент несколькими свойствами:

```
Name: Levi's Hooded Military Jacket
Description: Stand collar with drawstring hood
Brand: Levi Strauss
...
Price: 189.90
Category: Women's Jackets
```

Элементы также могут иметь динамические свойства, например данные о продажах и рейтинги пользователей, которые тоже содержат важную информацию о соответствии и, в конечном счете, релевантности. Значения свойств могут быть короткими строками, такими как названия продуктов, длинными фрагментами текста, например с описаниями или обзорами, числами, лексемами из дискретного набора, такими как названия брендов, или даже вложенными или иерархическими

сущностями, такими как разновидности продукта или категории. Это создает разнообразие характеристик и сигналов, измеренных на разных масштабах, которые нельзя сопоставлять непосредственно. Нам нужно найти способ, который поможет выявить корреляцию всех этих признаков с запросом и смешать полученные сигналы вместе, чтобы получить оценку релевантности.

Одно из простейших решений этой задачи заключается в объединении всех свойств в один большой текст и использовании базовых методов оценки для поиска по этому тексту. Это не самый бессмысленный подход, но он дает слабо выраженный сигнал, который непредсказуемо оценивает результаты поиска, опираясь на частоты термов и длину текста. Например, казалось бы, простой запрос *black dress shoes* (черные модельные туфли) может вернуть дикое сочетание платьев (*dress*), обуви (*shoes*), черных смокингов (*black tuxedos*) и других элементов, в описании которых присутствуют некоторые термы из запроса. Чтобы решить эту проблему, мы должны создать метод, сохраняющий главные признаки и сигналы и предоставляющий достаточно средств управления, чтобы выбрать самые точные и релевантные результаты.

4.4.1. Поиск по нескольким полям

Запрос можно рассматривать как описание желаемого результата поиска, в котором указывается, какими свойствами должен обладать элемент, чтобы считаться релевантным. Например, запрос *black levi strauss jeans* (черные джинсы Levi Strauss) явно запрашивает продукты типа *jeans* (джинсы) черного (*black*) цвета, выпускаемые под брендом *Levi Strauss*. Каждый элемент, в свою очередь, также представлен набором свойств, поэтому хорошего результата можно добиться, создавая документы с несколькими *полями*, каждое из которых соответствует свойству элемента, и отдельно проверять соответствие запросу каждого поля, чтобы получить несколько сигналов, а затем смешать их в окончательную оценку. Эту идею иллюстрирует рис. 4.8.

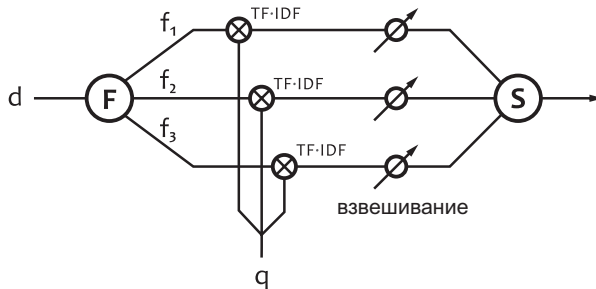


Рис. 4.8. Упрощенная схема оценки по нескольким полям. F обозначает деление документа на поля f_1, f_2, \dots, f_n . S — функция смешивания сигналов, вычисляющая окончательную оценку

Этот подход может давать достаточно хорошие результаты, но требует со всем вниманием отнестись к выбору четких и сбалансированных сигналов. Проблемы с четкостью сигналов могут возникать из-за несоответствия полей документа и понятий в пользовательском запросе, когда понятие, которое пользователь вкладывает в запрос, может оказаться разбросанным по нескольким полям. Например, пользователь может искать человека по имени и фамилии, ожидая, что они будут интерпретироваться как одна лексема, но документы могут хранить имена и фамилии в разных полях, создавая бессмысленные частичные совпадения. Проблема с балансом сигналов возникает из-за того, что каждое поле существует в своей собственной вселенной и нет общей шкалы для оценок сигналов (например, когда 0,00 означает нерелевантный, а 1,00 — релевантный).

ПРИМЕР 4.3

Рассмотрим проблему несбалансированности сигналов на примере каталога продавца модной одежды, содержащего тысячу товаров, в том числе джинсы и обувь:

Product 1

Name: Men's 514 Straight-Fit Jeans.

Description: Dark blue jeans. Blue denim fabric.

Brand: Levi Strauss

....

Product 1000

Name: Leather Oxfords.

Description: Elegant blue dress shoes.

Brand: Out Of The Blue

Можно ожидать, что для джинсов вполне обычным является присутствие термина *blue* (синий) в названии или описании, поэтому его значение IDF будет довольно низким. Например, если в упомянутой тысяче продуктов у нас имеется 500 синих джинсов, мы получим:

$$\text{idf}(\text{name:blue}) = 1 + \ln(1000 / 501) = 1,69. \quad (4.32)$$

В то же время бренд *Out Of The Blue*, под которым производится обувь синего цвета, может быть очень редким. Допустим, что в каталоге есть только один продукт этого бренда, и нет других брендов, содержащих слово *blue*, тогда значение IDF для термина *blue* в поле «brand» получится равным:

$$\text{idf}(\text{brand:blue}) = 1 + \ln(1000 / 2) = 7,21. \quad (4.33)$$

Теперь рассмотрим запрос *blue jeans*. Синие ботинки (*blue shoes*) получат очень высокую оценку релевантности для поля бренда и низкую — для поля описания, в то время как синие джинсы (*blue jeans*) получат относительно низкую оценку для поля описания, в котором присутствуют совпадения с обоими терминами из запроса, но имеющими низкое значение IDF. Комбинируя сигналы с помощью функции суммирования или выбора максимального значения, мы, скорее всего, получим список результатов с ботинками вверху, что явно не соответствует цели поиска. Причина в том, что оценки IDF зависят от распределения термов в поле, поэтому IDF для разных полей несопоставимы.

Одно из возможных решений проблемы несбалансированности сигналов заключается в определении весов вручную, как показано на рис. 4.8. В примере выше можно назначить низкий вес полю бренда, чтобы понизить мощность его сигнала и опустить обувь в конец списка. Такое решение может помочь в случаях, когда наблюдается последовательная разность в уровнях или важности сигналов (например, когда соответствие названий важнее, чем соответствие описаний), но это очень хрупкое решение. К проблеме регулировки сигналов нужно подходить более системно.

4.4.2. Проектирование и регулировка сигналов

Поиск по нескольким полям имеет два аспекта, дополняющих друг друга, которые, как правило, должны учитываться одновременно, — проектирование сигналов и регулировка сигналов. Целью проектирования сигналов является создание ясных и четких сигналов, а целью регулировки сигналов является правильное их смешивание для получения окончательного результата. Иногда одну и ту же задачу релевантности можно решить по-разному, — либо путем настройки функции смешивания, либо путем проектирования более точных сигналов. В поиске по нескольким полям различаются следующие типы отношений между полями и целью поиска [Gormley and Tong, 2015]:

- *Один сильный сигнал.* Случай, когда пользователь ищет определенное свойство, которое в идеале должно соответствовать одному из полей и производить один сильный сигнал. Сигналы из разных полей не дополняют, а скорее конкурируют друг с другом. Например, пользователь, который ищет бренд *Out Of The Blue*, скорее всего, будет сосредоточен на поле бренда и не считает синий (*blue*) цвет релевантным.

- *Сильный средний сигнал.* Иногда предпочтительнее выбирать средний сигнал, а не самый сильный, если отдельные сигналы сбалансированы и связаны с различными аспектами одной и той же цели поиска. Например, размер и цвет элемента могут быть одинаково важны.
- *Фрагментированные признаки и сигналы.* Получить четкий сигнал, оценивая отдельные поля, можно, только если запрос и поле согласованы и созвучны. Однако может случиться так, что поля содержат фрагментированные обрывки информации, из-за чего получающиеся сигналы не коррелируют с релевантностью. Такие фрагменты можно объединить и получить более сильный сигнал.

Эти три типа отношений тесно связаны между собой. Давайте пройдемся по списку и подробно обсудим методы проектирования и регулировки сигналов для каждого случая.

4.4.2.1. Один сильный сигнал

Проблема с брендом *Out Of The Blue*, о которой говорилось в предыдущем разделе, возникает из-за неточной обработки сигнала бренда. Мы отметили, что одним из возможных решений является регулировка веса сигнала, однако есть альтернативный путь — уточнение сигнала бренда с целью сделать его менее двусмысленным. Можно утверждать, что название бренда — это понятие, которое не должно разбиваться на отдельные слова, поэтому есть смысл заменить оценку $TF \times IDF$ оценкой на основе биграмм. Это усилит сигнал, если в запросе присутствует узнаваемая часть названия бренда. Также для смешивания сигналов можно использовать функцию выбора максимального значения, чтобы выбрать самый сильный сигнал. В этом случае мы приходим к конвейеру оценки, представленному на рис. 4.9.

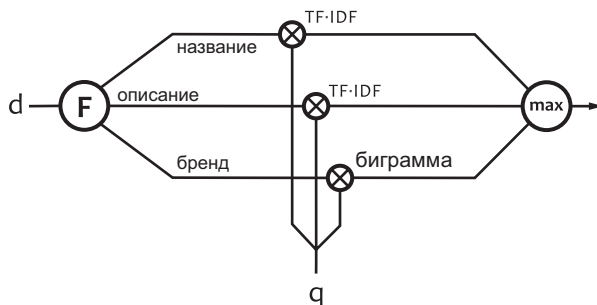


Рис. 4.9. Пример конвейера смешивания сигналов, выделяющего самый сильный сигнал

Результат поиска с использованием конвейера 4.9 показан на рис. 4.10. Товары, продаваемые под брендом *Out Of The Blue*, превосходят другие товары, только если название бренда четко сформулировано в запросе; в противном случае приоритет получают товары с релевантными описаниями и названиями. Обратите внимание, что нам может потребоваться внести дополнительные изменения, чтобы оценка на основе биграмм работала правильно, например отфильтровать стоп-слова.

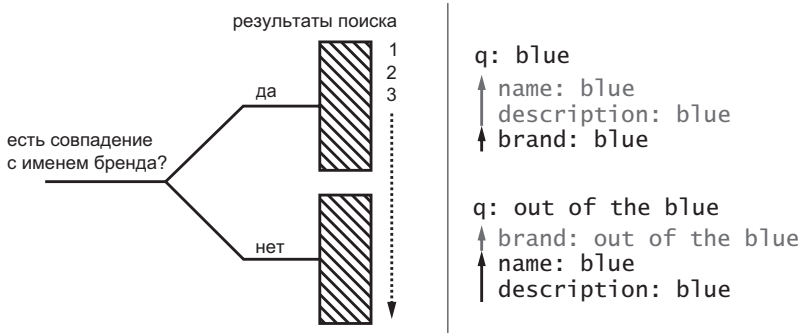


Рис. 4.10. Структура результатов поиска с использованием стратегии выбора сильного сигнала

Стратегия смешивания сигналов с использованием функции выбора максимального значения — мощный и популярный подход к поиску по нескольким полям, а использование n -грамм — эффективный метод фокусировки сигнала. Однако создание сфокусированных сигналов — сложный процесс, который не ограничивается применением n -грамм. Нет ничего необычного в том, чтобы генерировать несколько сигналов на основе одного и того же поля, например, с использованием униграмм и биграмм, которые нужно смешать, или попросить мерчандайзеров снабдить элементы совершенно новым свойством, которое поможет получить более сфокусированный сигнал, чем уже имеющиеся признаки.

Структура результатов поиска на рис. 4.10 довольно проста, и в действительности можно запрограммировать более сложное поведение, ослабив функцию смешивания сигналов. Одно из возможных решений — смешивание слабых и сильных сигналов некоторым управляемым способом. Это можно выразить с помощью следующей формулы:

$$s = s_m + \alpha \sum_{i \neq m} s_i, \quad (4.34)$$

где s_m — максимальный (самый сильный) сигнал, $0 \leq \alpha \leq 1$ — параметр, контролирующий вес всех остальных сигналов s_i в смеси, а s — окончательная оценка. Легко

видеть, что формула 4.34 представляет целый спектр функций оценки, начиная с выбора самого сильного сигнала, когда α равно нулю, и заканчивая усреднением сигнала, когда α равно единице. Такой подход позволяет получить структуру результатов с более чем двумя уровнями релевантности. Например, основной приоритет можно отдать продуктам с названиями, соответствующими запросу, и придать дополнительные очки за соответствие в поле бренда, как показано на рис. 4.11. Реализовать это можно с помощью функции 4.34, настроив веса так, чтобы усилить сигнал названия продукта и поднять соответствующие элементы ближе к началу списка с результатами поиска. Соответствие бренда будет вторым по силе сигналом в смеси, поэтому элементы во внутренних уровнях, созданные сопоставлением имен, будут ранжироваться по бренду. Конвейер смешивания сигналов, реализующий эту стратегию, показан на рисунке 4.12.

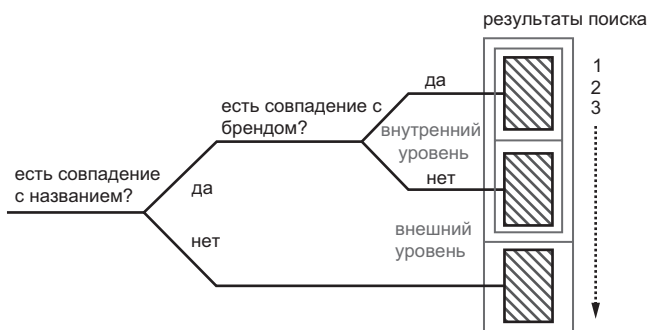


Рис. 4.11. Структура результатов поиска для взвешенного смешивания сигналов

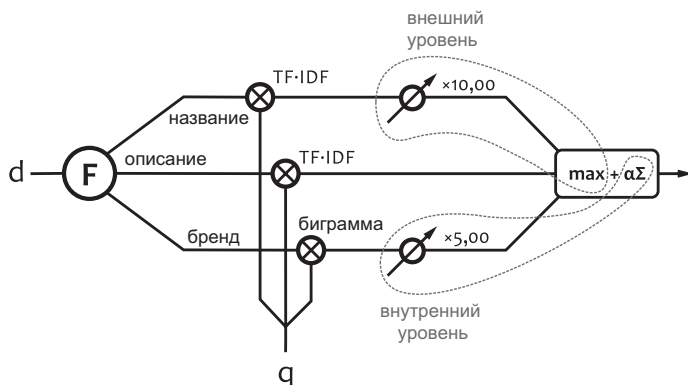


Рис. 4.12. Пример конвейера для взвешенного смешивания сигналов

4.4.2.2. Сильный средний сигнал

Если усилить вторичные сигналы в формуле 4.34, получится решение, ранжирующее элементы по среднему сигналу. Сумма сигналов тоже дает своего рода многоуровневые результаты поиска, хотя в этом случае количество уровней определяется количеством критериев поиска, которым соответствует элемент, а не количеством полей. Поисковый запрос может содержать несколько критериев, соответствующих разным свойствам элемента и, следовательно, полям документа. Суммируя сигналы из разных полей, мы получим ранжирование, которое помещает элементы со многими созвучными свойствами в верхнюю часть списка результатов, за которыми следуют элементы с меньшим количеством свойств, созвучных запросу. Мы уже видели, что этот способ подвержен проблеме дисбаланса сигналов, но он может быть подходящим решением для смешивания группы скорректированных сигналов. Обычно этот метод используется с сигналами, полученными из разных версий одного и того же свойства. Например, можно применить стемминг к определенному полю, чтобы улучшить полноту, и одновременно использовать оригинальную версию того же поля без стемминга, для повышения веса точных совпадений. Следующие описания продукта содержат слова *fashion* (мода) и *fashionable* (модно):

d_1 : new popular fashion brand
 d_2 : stylish and fashionable look

которые стеммером Портера сокращаются до термина *fashion* (мода). Как результат, оба описания получают одинаковые оценки $TF \times IDF$ для запроса *fashionable* (модно). Оценку релевантности для второго документа, в котором есть слово, точно соответствующее запросу, можно увеличить, добавив оценку за совпадение без стемминга и сложив оценки за совпадение со стеммингом и без стемминга. Суммирование сигналов в данном случае вполне обоснованно, потому что чем больше совпадений, тем лучше результат.

4.4.2.3. Фрагментированные признаки и сигналы

Наконец, обсудим случай фрагментированных признаков, когда требуется применение методов проектирования сигналов. Проблема фрагментированных признаков и сигналов возникает из-за того, что стандартный поиск по нескольким полям, о котором говорилось выше, оценивает каждое поле независимо. Хотя на первый взгляд это не всегда очевидно, но иногда отдельные поля достаточно хорошо коррелируют с запросом, чтобы получить сильный средний сигнал, но общий охват условий запроса всеми полями остается низким. Рассмотрим простой пример с двумя описаниями продуктов [Turnbull and Berryman, 2016]:

Product 1

Name: Polo

Brand: Polo

Product 2

Name: Polo

Brand: Lacoste

Очевидно, что второй продукт более релевантен запросу *Lacoste Polo*, но вычисление оценок $TF \times IDF$ дает иной результат. Напомню, что практическая формула 4.29 оценки $TF \times IDF$ для полей с одним термом сводится к следующему уравнению:

$$\frac{\text{координационный фактор запроса}}{\text{норма запроса} \times \text{норма поля}} \times \text{tf}(\text{term, field}) \times \text{idf}_{\text{field}}^2(\text{term}). \quad (4.35)$$

Координационный фактор запроса (query coordination factor) равен 0,50 для всех четырех полей, потому что совпадение в них обнаруживается только с одним из термов в запросе (*Polo* или *Lacoste*). Норма запроса (query norm), норма поля (field norm) и значение TF также одинаковы для всех термов и полей. Значения IDF для *Polo* и *Lacoste* одинаковы для поля бренда, но отличаются от IDF *Polo* для поля с названием. Следовательно, поля с названием и брендом (попарно) получат одинаковые оценки в обоих документах, и итоговые оценки документов также будут равны, какая бы функция ни использовалась для смешивания сигналов, суммирования или выбора максимального значения. Основная причина в том, что все поля соответствуют ровно одному терму запроса (*Polo* или *Lacoste*) и дают одинаково сильные сигналы, и при этом не учитывается тот факт, что второй документ включает два терма, а первый — только один. Эта проблема может возникнуть во многих случаях, когда разные аспекты одного и того же логического свойства моделируются как разные поля: имя человека может быть разбито на имя и фамилию, адрес доставки может быть разделен на название улицы, название города и страны, и так далее. Фрагментированные сигналы могут приводить к расстраивающим результатам поиска — документ, идеально соответствующий запросу, может присутствовать в списке результатов поиска, но иметь удивительно низкий ранг.

Одно из решений проблемы заключается в том, чтобы объединить несколько похожих полей в одно и таким способом устранить проблему фрагментированных или несбалансированных сигналов. Это действенный и практичный метод, способный повысить релевантность. Но он имеет свой недостаток — результирующий сигнал может получиться слишком размытым, из-за того что вторичные объекты становятся такими же важными, как первичные. Например, описание платья может содержать фразу *wear with any shoes* (можно носить с любой обувью), способную поднять платье вверх в результатах поиска обуви, если признак *тип продукта* не взвешивается должным образом.

Альтернативное решение проблемы фрагментированных признаков основано на наблюдении, что запросы с одним термом не подвержены фрагментации, описанной выше. Кроме того, каждый терм в запросе можно рассматривать как отдельный критерий, добавляемый пользователем для сужения поиска, поэтому иногда разумно получить оценку документа для каждого терма в запросе независимо, создав сигнал, указывающий степень соответствия этому конкретному критерию, а затем смешать все сигналы, чтобы получить окончательную оценку. Этот подход назвали *оценкой по термам* (term-centric scoring), чтобы отличить его от подхода оценки по полям (field-centric), который мы использовали выше. Конвейер оценки по термам можно рассматривать как группу конвейеров оценки по полям, выполняемых для получения сигналов по каждому терму запроса, которые затем смешиваются в конечный результат, как показано на рис. 4.13. Сигналы из разных конвейеров суммируются, потому что чем больше совпадающих термов, тем лучше (стратегия среднего сигнала); между тем к сигналам из разных полей в пределах конвейера для одного терма можно применить стратегию одного сильного сигнала.

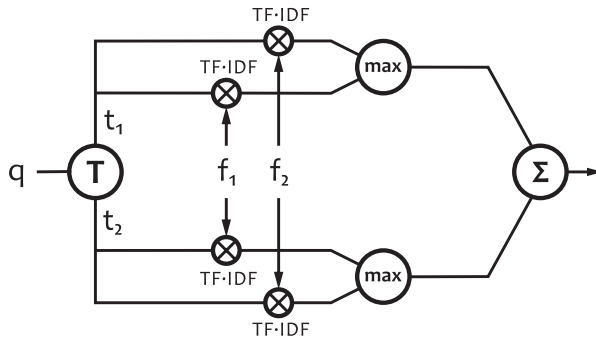


Рис. 4.13. Конвейер оценки по термам. t_1 и t_2 — термы в запросе; f_1 и f_2 — поля в документах

Вернемся к нашему примеру с брендами *Polo* и *Lacoste*: обратите внимание, что подход оценки по термам даст более значимые результаты для запроса *Lacoste Polo*. Первый документ получит высокую оценку за совпадение с термом *Polo* в обоих полях и нулевую оценку за отсутствие терма *Lacoste*, поэтому общая оценка получится равной

$$\max\{tf - igf_{name}(Polo), tf - igf_{brand}(Polo)\} + \max\{0, 0\}, \quad (4.36)$$

где оценка $TF \times IDF$ терма *Polo* для поля с названием составит 1,00 и для поля бренда составит 0,35 (из-за различий в величине IDF). В то же время второй документ получит высокую оценку за совпадение с обоими термами в запросе:

$$\max\{\text{tf} - \text{igf}_{\text{name}}(\text{Polo}), 0\} + \max\{0, \text{tf} - \text{igf}_{\text{brand}}(\text{Lacoste})\}, \quad (4.37)$$

где оценка TF×IDF термина *Lacoste* для поля бренда составит 1,00. Этот результат выглядит более релевантным, чем тот, что получается при использовании подхода к оценке по полям.

4.4.3. Проектирование конвейера смешивания сигналов

Мы видели, что структура результатов поиска определяется архитектурой конвейера смешивания сигналов. Этот процесс можно перевернуть и попытаться получить конвейер по известной структуре результатов. Такая постановка задачи имеет большое практическое значение, поскольку позволяет проектировать признаки и функции оценки по описанию желаемого результата. Этот подход тесно связан с проектированием релевантности и управлением продвижением, поскольку описание может включать предметные знания о критериях релевантности и бизнес-целях. Программная система может предоставлять интерфейс, облегчающий описание желаемых результатов поиска и проектирование конвейеров смешивания сигналов, а также экспериментальную оценку.

Рассмотрим пример относительно сложного описания результатов поиска, чтобы продемонстрировать сквозное проектирование сигналов и функций оценки с использованием текстовых и нетекстовых признаков. Возьмем за основу ретейлера, торгующего модной одеждой, который решил создать службу онлайн-поиска. Также в целях упрощения предположим, что пользователь выполняет поиск в определенной категории (как достичь высокой точности при поиске по нескольким категориям, мы обсудим в одном из следующих разделов). Отправной точкой для нас послужит описание, представленное на рис. 4.14, которое кодирует следующие бизнес-правила:

- Если пользователь ищет определенный товар по названию или идентификатору, соответствующий товар должен находиться в верхней части результатов поиска.
- Если пользователь ищет определенный бренд, товары этого бренда должны иметь приоритет и дополнительно сортироваться по новизне и рейтингу, определяемому клиентами, чтобы поднять вверх новые товары и товары с высоким рейтингом.
- В противном случае результаты должны ранжироваться в соответствии со средней релевантностью описаний товаров и других полей.

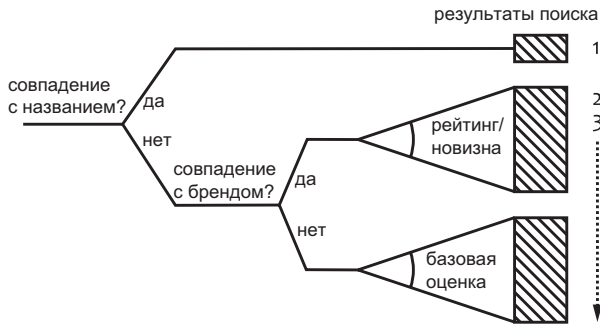


Рис. 4.14. Пример описания результатов поиска, которое можно использовать для проектирования конвейера смешивания сигналов

В описании упоминается 5 разных сигналов, которые требуется учесть: точное совпадение с названием товара или его идентификатором, точное совпадение с брендом, новизна товара, рейтинг товара и базовая средняя оценка. Точные сигналы можно получить путем сопоставления соответствующих полей названия и бренда с n -граммами. Достаточно хорошую точность можно получить с использованием биграмм, а еще более высокую точность — с использованием триграмм или логического сопоставления фраз. Обратите внимание, что для n -грамм можно даже отключить вычисление $TF \times IDF$ и просто подсчитать количество совпавших n -грамм, потому что нас интересует бинарный результат (есть совпадение с названием товара или нет?), а не непрерывные оценки релевантности. Предположим, что новизна и рейтинг доступны как числовые свойства товара. Новизна, например, может измеряться как количество дней с момента появления товара в магазине, а средний покупательский рейтинг может быть вещественным числом в шкале от 1 до 5. Базовую оценку можно рассчитать с использованием одного из методов проектирования сигналов, о которых говорилось выше. Проще всего было бы объединить все свойства товара в одно поле и вычислить для него оценку $TF \times IDF$.

Начнем конструирование конвейера смешивания сигналов с определения базовой оценки, и затем будем добавлять новые уровни релевантности, в соответствии с описанием результатов поиска. Каждый новый уровень будет усиливать соответствующий сигнал, чтобы все товары, обладающие свойствами, необходимыми для этого сигнала, поднимались в результатах поиска выше нижних уровней.

Первый уровень, который мы надстроим над базовой оценкой, — точное совпадение с названием бренда. Для достижения желаемой вторичной сортировки мы должны усилить сигнал бренда, назначив ему повышающий коэффициент, а также смешать его с признаками рейтинга и новизны, как показано на рис. 4.15. Очевидно,

что исходные значения рейтинга и новизны желательно масштабировать, чтобы преобразовать их в значимые факторы оценки; сделать это можно по-разному. Исходный рейтинг клиентов по шкале от 1 до 5 может дать слишком агрессивный коэффициент усиления, поэтому его лучше сгладить, например, вычислив квадратный корень или логарифм, и тем самым уменьшить разрыв между товарами с низкими и высокими рейтингами. Например, величина сигнала совпадения с названием бренда, усиленного исходным рейтингом 5.0, будет в два раза выше величины сигнала, усиленного рейтингом 2.5; однако, взяв квадратный корень рейтинга, мы уменьшаем эту разницу до 1,41 раза.

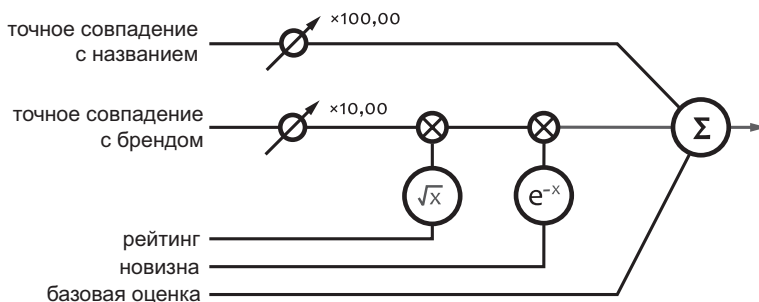


Рис. 4.15. Конвейер смешивания сигналов, соответствующий описанию на рис. 4.14

Значение новизны следует преобразовать в коэффициент, постепенно уменьшающийся с увеличением времени, прошедшего с момента выпуска товара. Это можно сделать с помощью линейной функции экспоненциального спада или гауссова спада. Например, вполне разумно уменьшать этот коэффициент на 10 % за каждые 30 дней, что даст нам экспоненциально затухающую функцию:

$$\text{newness factor} = \exp(-\alpha x), \quad (4.38)$$

где x — значение новизны в днях, а параметр α определяется из следующего уравнения (согласно которому отношение между любыми двумя факторами, отстоящими друг от друга на 30 дней, равно 0,9):

$$\exp(-30 \cdot \alpha) = 0,9. \quad (4.39)$$

Комбинируя сигнал совпадения с названием бренда и коэффициентами рейтинга и новизны, получаем сигнал, соответствующий второму уровню. Наконец, верхний уровень формируется на основе сигнала совпадения с названием товара, смешанного с большим постоянным коэффициентом, который поднимает товары с совпадающими названиями в самый верх списка с результатами поиска.

4.5. Семантический анализ

Методы поиска, рассматривающиеся до сих пор, основаны на сопоставлении лексем. Несмотря на использование таких методов, как стемминг, учитывающих особенности естественного языка, мы фактически свели проблему поиска к механическому сравнительному лексем. Этот подход, иногда называемый *синтаксическим поиском*, хорошо работает на практике и используется как основной метод в большинстве поисковых систем. Однако синтаксический поиск имеет ограниченные возможности моделирования особенностей естественного языка, которые выходят за рамки отдельных термов. Значение слов в естественном языке часто зависит от контекста, созданного предыдущими и последующими словами и предложениями, и существует несколько типов таких зависимостей. Большинство из них относится к одной из двух категорий:

ПОЛИСЕМИЯ. Полисемия (многозначность) — это ситуация, когда слово имеет несколько значений. Например, слово *дерево* может обозначать материал или растение. Полисемия представляет серьезную проблему для релевантности, потому что пользователь может иметь в виду одно значение слова (например, изделия из дерева), а поисковая система будет возвращать документы, где то же слово используется в другом значении (например, с описанием оборудования для использования в лесном хозяйстве). Мы уже сталкивались с проблемой сложных понятий, выражаемых фразами, которые необходимо интерпретировать как единую лексему, например, *dress shoes* (модельные туфли). Эту проблему можно рассматривать как частный случай полисемии, потому что значение отдельных слов зависит от контекста — например, когда значение слова *dress* (платье) зависит от того, следует ли за ним слово *shoes* (обувь). Частным случаем полисемии, которая очень распространена в сфере продвижения товаров, является использование словаря слов, встречающихся в названиях брендов и товаров. Примером может служить название бренда *Blue*, которое неотличимо от слова *blue* (синий) в таких запросах, как *blue jeans* (синие джинсы).

СИНОНИМИЯ. Слова считаются синонимами, если выражают почти то же значение в данном контексте, как, например, *конфеты* и *леденцы*. Синонимия также является серьезной проблемой для релевантности поиска, поскольку простой синтаксический поиск не в состоянии выбрать релевантные документы, которые не содержат данного термина из запроса, но содержат его синонимы. Например, документы, содержащие слово *леденцы*, почти наверняка будут релевантны запросу *конфеты*, но все методы поиска, о которых говорилось до сих пор, не способны справиться с этой зависимостью.

Некоторые аспекты проблемы полисемии можно решить с применением *n*-грамм и более продвинутых методов сопоставления фраз, которые мы обсудим в следу-

ющих разделах, хотя бы отчасти. Аналогично стемминг можно рассматривать как метод решения проблемы синонимии в простейшей ее форме, путем приведения близких термов к одному корню. Эти методы, однако, не учитывают значений слов и взаимосвязей между ними, что необходимо для решения проблем полисемии и синонимии. Чтобы решить эти проблемы, мы должны разработать новые методы, ориентированные на анализ контекста и смысла, а не отдельных лексем. Этот подход известен как *семантический поиск*, названный в честь раздела лингвистики, изучающего значения слов и отношения между словами и фразами. Некоторые семантические методы полностью независимы от синтаксического поиска и конкурируют с ним, но при этом многие из них можно использовать для расширения синтаксического поиска.

Полисемию и синонимию можно рассматривать как задачу выявления скрытых отношений между словами или, наоборот, как задачу выявления логических понятий, материализованных в словах. С последней точки зрения полисемия относится к случаю отображения двух разных понятий в одном слове, тогда как синонимия — к противоположному случаю, когда два разных слова отражают одно логическое понятие. Соответственно, задачу семантического анализа и поиска можно рассматривать как задачу выявления правильных понятий и связей между словами и понятиями. Такой способ мышления часто называют *концептуальным (ассоциативным) поиском* [Giunchiglia et al., 2009; Hughes, 2015]. Этот термин подчеркивает тот факт, что понятия являются не просто статистическими отношениями между словами, а логическими сущностями, которые можно определять, руководствуясь знаниями о предметной области и другими соображениями.

В оставшейся части этого раздела мы рассмотрим методы семантического поиска и анализа, которые помогут решить проблемы полисемии и синонимии, причем некоторые из этих методов тесно связаны с выработкой рекомендаций. Также для разнообразия заменим в примерах одежду бакалейными продуктами.

4.5.1. Синонимы и иерархии

Самым простым решением проблемы синонимии является поддержка *тезауруса*, составляемого вручную, то есть каталога слов и их синонимов. После создания тезаурус можно использовать для преобразования документов и запросов способом, похожим на стемминг. Например, можно определить следующий набор синонимов:

candy, sweet, confection

Наша цель — сделать эти термы идентичными с точки зрения запроса, чтобы документы, содержащие слова *sweet* (леденцы) и *confection* (сладости), извлекались по

запросу *candy* (конфеты) и наоборот. Это можно сделать несколькими способами, каждый из которых имеет свои достоинства и недостатки.

Первый подход — *сужение* (contraction). Один из синонимов в списке назначается основным, и все другие синонимы заменяются им. Например, все слова *sweet* (леденцы) и *confection* (сладости) в документах и в запросах можно заменить словом *candy* (конфеты). Обратите внимание, что основным синоним не обязательно должен быть реальным словом; это может быть специальная лексема, никогда не появляющаяся во входных текстах, но используемая как внутреннее представление группы синонимов. Таким образом, сужение работает точно так же, как стемминг. Сужение явно преследует цель сделать все синонимы идентичными с точки зрения запроса, и главный его недостаток состоит в том то, что этот подход приводит все синонимы к основному синониму, делая часто используемые синонимы неотличимыми от редко используемых. Это может негативно отразиться на оценке $TF \times IDF$.

Альтернативой сужению является *расширение* (expansion). Согласно стратегии расширения, каждое вхождение синонима заменяется полным списком синонимов:

`best candy shop → [best] [candy] [sweet] [confection] [shop]`

Расширение может применяться к документам или запросам, но не к обоим сразу. Расширение на стороне документа может оказать такое же негативное влияние на оценку $TF \times IDF$, как и сужение, а также увеличить размер документов. Расширение на стороне запроса сохраняет правильную статистику IDF , но делает запрос более сложным с вычислительной точки зрения.

Прием расширения имеет одно очень важное применение, выходящее далеко за пределы простой обработки синонимов. Хотя синонимы определяются как слова, имеющие примерно одинаковое значение, часто бывает так, что одно слово представляет более широкое логическое понятие, чем другое. Связь между такими синонимами становится асимметричной в том смысле, что более широкое понятие можно рассматривать как синоним более узкого, но не наоборот. Например, слово *торт* можно использовать как синоним понятия *творожный торт*, но в некоторых контекстах было бы неправильно заменить слово *торт* словосочетанием *творожный торт*. Следовательно, было бы полезно проработать процесс расширения и заменить простые списки синонимов правилами, допускающими замену *творожного торта* словом *торт*, но запрещающими обратное расширение. Этот тип расширения называется *жанровым расширением*.

Продолжая развивать эту идею, можно построить иерархию термов, описывающую вложенные классы понятий, как показано на рис. 4.16. На каждом уровне иерархии термы расширяют своих предшественников и благодаря этому могут обнаруживаться запросами, содержащими более общие термы. Например, элемент,

содержащий терм *фруктовый пирог*, будет включен в список результатов поиска по запросам с терминами *торт* и *выпечка*.

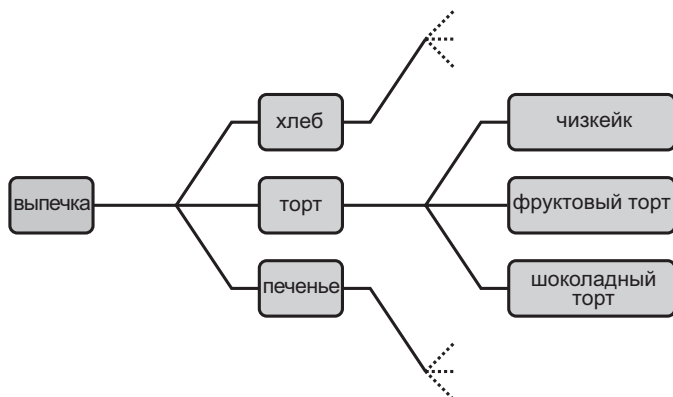


Рис. 4.16. Пример иерархии категорий бакалейных продуктов

Подобно простому расширению, жанровое расширение можно применять и к документам, и к запросам, но методология сильно отличается из-за асимметрии правил расширения. Применяя жанровое расширение к документам, мы делаем запросы с общими понятиями соответствующими документам с конкретными понятиями, но не наоборот. Например, запрос *торт* вернет *творожные торты* и *фруктовые пироги*. Применяя расширение к запросам, мы делаем запросы с конкретными понятиями, соответствующими документам с широкими понятиями, то есть запрос *творожный торт* вернет *торты*. Применяя жанровое расширение к документам, мы также искусственно уменьшаем IDF широких понятий, потому что эти термины копируются в другие документы. Это не обязательно плохо, потому что конкретные понятия будут оцениваться выше, чем широкие, что обычно правильно с точки зрения релевантности.

Сужение и расширение — очень мощные методы моделирования семантических сетей. Они подсказывают, как можно использовать известные семантические отношения в обработке запросов. В то же время они не содержат указаний, как вывести эти отношения. Одно из возможных решений — создание списков синонимов вручную. Это часто делается в приложениях поиска для продвижения, поскольку позволяет мерчандайзерам использовать синонимы как элемент управления, способный выражать определенные бизнес-правила и знания предметной области. С другой стороны, ручное сопровождение тезауруса в приложениях поиска с динамическим контентом может оказаться сложной задачей. Кроме того, порой трудно выявить некоторые виды семантических зависимостей без применения

методов машинного обучения. Например, имя известного спортсмена может ассоциироваться с определенным видом спорта, спортивным инвентарем или брендом, который продвигает этот спортсмен. Нашей следующей целью станет разработка методов, способных автоматически изучать тезаурус.

4.5.2. Векторные представления слов

Модель векторного пространства утверждает, что документ или запрос можно представить в виде вектора в линейном пространстве термов. Рассматривая проблемы полисемии и синонимии, мы узнали, что термы могут быть многозначными и избыточными, и может случиться так, что представление документа, использующее термы в качестве измерений, окажется не самым лучшим или, по крайней мере, будет иметь некоторые недостатки. И действительно, мы уже выяснили, что полисемию и синонимию можно рассматривать как проблемы несоответствия слов и понятий, предполагающие, что слова имеют запутанные представления, скрывающие семантические отношения.

Мы можем попытаться найти лучшее представление, изменив основу для пространства документов. С концептуальной точки зрения хотелось бы иметь возможность представлять документы и запросы в виде векторов вещественных чисел, чтобы вычислять ранжирующие оценки как простые скалярные произведения между представлениями запросов и документов:

$$\begin{aligned} q &\rightarrow \mathbf{p} \\ d &\rightarrow \mathbf{v} \end{aligned} \tag{4.40}$$

$$\text{score}(q, d) = \mathbf{p} \cdot \mathbf{v} = \sum_{i=1}^k p_i v_i,$$

где \mathbf{p} и \mathbf{v} являются векторными представлениями запроса и документа, а k — размерность векторного представления. Обратите внимание, что отдельные слова также могут быть представлены такими же векторами, потому что каждое слово можно рассматривать как документ из одного слова. Этот подход часто называют *векторным представлением слов* (word embedding). Число измерений k обычно мало в сравнении с числом различных слов, поэтому модель векторного представления пространства высокой размерности, в котором каждое слово имеет собственное измерение, как бы *отражается* в пространство малой размерности. Векторное представление слов имеет много применений в зависимости от того, как устроены и используются векторные представления. С точки зрения служб поиска основной интерес представляют две возможности. Во-первых, векторные представления можно использовать для фактической обработки запросов, чтобы оценки документов для данного запроса вычислялись как скалярные произведе-

ния представлений. Во-вторых, векторные представления отдельных слов можно проанализировать для создания тезауруса, то есть для сопоставления слов с их синонимами или связанными с ними словами.

Ключевой проблемой векторных представлений слов является создание новых векторных представлений. С концептуальной точки зрения нам нужно векторное пространство, сохраняющее семантические отношения: слова и документы с похожими или связанными семантическими значениями должны располагаться, или лежать, на одной линии, а документы с разными семантическими значениями не должны быть коллинеарными, даже если содержат одни и те же (многозначные) слова. Сумев построить такое пространство, мы смогли бы преодолеть ограничения векторной модели пространства. Во-первых, мы смогли бы находить документы, соответствующие запросу, даже если бы они не имели общих термов, то есть решили бы проблему синонимии. Во-вторых, мы смогли бы исключить семантически нерелевантные документы, даже если бы они номинально содержали термы из запроса. Интуитивно можно предположить, что семантическое пространство можно построить путем анализа словосочетаний либо глобально в документе, либо локально в предложении, и определив группы связанных слов. Затем на основе этих групп можно определить размеры семантического пространства, а векторные представления отдельных документов и слов, соответственно, в терминах сходства с группами. Оказывается, эту простую идею очень сложно реализовать, и существует большое количество методов, которые используют очень разные математические методы. Далее в этом разделе мы подробно обсудим несколько важных подходов и конкретных моделей.

Наконец, немного о терминологии. Векторное представление слов — относительно новый термин, и многие методы семантического анализа, включая латентно-семантический анализ и вероятностное тематическое моделирование, описываемые далее, разрабатывались не специально для векторного представления слов (в смысле уравнения 4.40), а для других целей и на основе других соображений. Большинство из этих методов являются очень мощными и универсальными статистическими методами, используемыми в широком диапазоне применений, от обработки естественного языка до эволюционной биологии. Однако эти методы также можно рассматривать как методы создания векторных представлений слов. В этом разделе мы будем рассматривать в основном векторные представления слов, потому что они дают удобный способ связать друг с другом разные семантические методы, по крайней мере в контексте поиска для продвижения. Но имейте в виду, что это только одна возможная перспектива; методы семантического анализа не ограничиваются векторными представлениями слов и поиском, так же как применение векторных представлений слов не ограничивается приложениями поиска. Даже в сфере алгоритмического маркетинга методы семантического анализа могут применяться во многих областях, включая автоматизированную оценку продукта, выработку рекомендаций и поиск изображений.

4.5.3. Латентно-семантический анализ

Один из возможных подходов к конструированию семантического пространства — выполнить анализ словарных представлений документов и выяснить, какие термины чаще всего встречаются вместе в одном документе. Интуитивно понятно, что термины, часто встречающиеся вместе, могут быть синонимами, соответствующими одному логическому *понятию*. То есть анализ совместных вхождений может выявить понятия, которые не наблюдаются в документах явно как термины, а существуют на семантическом уровне. Эти понятия называются *скрытыми (латентными) понятиями*.

Начнем с набора документов с единственным полем, содержащим текстовое описание товара. Первым шагом подготовим матрицу с частотами термов для всех терминов t_i и документов d_j :

$$\mathbf{X} = \begin{matrix} & \begin{matrix} d_1 & d_2 & \dots & d_n \end{matrix} \\ \begin{matrix} t_1 \\ t_2 \\ \vdots \\ t_m \end{matrix} & \begin{bmatrix} \text{tf}(t_1, d_1) & \text{tf}(t_1, d_2) & \dots & \text{tf}(t_1, d_n) \\ \text{tf}(t_2, d_1) & \text{tf}(t_2, d_2) & \dots & \text{tf}(t_2, d_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{tf}(t_m, d_1) & \text{tf}(t_m, d_2) & \dots & \text{tf}(t_m, d_n) \end{bmatrix} \end{matrix}, \quad (4.41)$$

где n — количество документов, а m — общее количество разных термов в коллекции. Эта матрица, известная как *матрица термов/документ*, содержит представления документов в пространстве терминов. Мы можем определить сходство между документами, вычислив скалярное произведение между соответствующими столбцами, и сходство между термами, вычислив скалярное произведение между соответствующими строками. Сходство термов, вычисленное таким способом, уже дает некоторые подсказки о семантических связях между термами, в том смысле что термы, которые часто встречаются вместе, могут быть связаны с одним и тем же понятием. Эта метрика, однако, может оказаться очень неточной, поэтому нам нужен более надежный статистический метод.

Напомним, что парадигма векторного представления слов предполагает представление каждого документа в виде k -мерного вектора. Это представление также можно записать в виде матрицы. Определим его как матрицу \mathbf{V}_k с размерами $n \times k$, в которой каждая строка соответствует документу, а каждый столбец — семантическому измерению. Метод латентно-семантического анализа (Latent Semantic Analysis, LSA) создает эту матрицу, основываясь на эвристическом рассуждении, что должна иметься возможность восстановить приближительную матрицу термов/документ \mathbf{X} из \mathbf{V}_k с помощью линейного преобразования [Deerwester et al., 1990]. Другими словами, должна быть возможность вычислить $m \times k$ матрицу \mathbf{L}_k , такую, что

$$\mathbf{X} \approx \mathbf{L}_k \cdot \mathbf{V}_k^T. \quad (4.42)$$

Более конкретно, матрицы \mathbf{L}_k и \mathbf{V}_k должны быть определены так, чтобы минимизировать ошибку восстановления. Если в качестве меры использовать среднеквадратичную ошибку, тогда мы приходим к следующей задаче оптимизации:

$$\min_{\mathbf{L}_k, \mathbf{V}_k} \|\mathbf{X} - \mathbf{L}_k \cdot \mathbf{V}_k^T\|. \quad (4.43)$$

Напомним, что эта задача решается методом сингулярного разложения (SVD, см. главу 2). Кроме того, очень важно, что столбцы в матрице, создаваемой алгоритмом SVD, являются ортонормальными, то есть все k размерностей (столбцы матрицы \mathbf{V}_k) будут ортогональны друг другу. По сути, это означает, что исходные векторы модели векторного пространства (строки матрицы \mathbf{X}) будут декоррелированы путем свертывания сильно коррелированных векторов в один главный вектор. Это соответствует нашему интуитивному ожиданию — термы, часто встречающиеся вместе, соответствуют сильно коррелированным компонентам исходных векторов термов (строки матрицы \mathbf{X}), поэтому декорреляция почти наверняка объединит термы, встречающиеся вместе (потенциально синонимы), в один вектор понятия.

Опишем этот процесс более формально. Сначала рассмотрим случай полного сингулярного разложения, в котором число размерностей понятия k не ограничено. Алгоритм SVD разбивает матрицу на три множителя:

$$\begin{aligned} \mathbf{X} &= \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \\ &= \begin{bmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_r \\ | & & | \end{bmatrix}_{m \times r} \begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r \end{bmatrix}_{r \times r} \begin{bmatrix} - & \mathbf{v}_1 & - \\ & \vdots & \\ - & \mathbf{v}_n & - \end{bmatrix}_{n \times r}^T, \end{aligned} \quad (4.44)$$

где r — ранг матрицы термов/документ \mathbf{X} . Давайте внимательно рассмотрим это разложение, чтобы понять, как оно может пригодиться в семантическом анализе и поиске.

Столбцы матрицы \mathbf{U} можно интерпретировать как новую основу для пространства документов. Каждый столбец можно рассматривать как скрытое понятие, которое может включать несколько коррелированных термов, то есть термов, которые часто встречаются вместе в одном документе. Каждая строка матрицы \mathbf{U} соответствует терму, соответственно элемент u_{ij} определяет важность или вклад термина t_i в понятие \mathbf{u}_j . Может случиться так, что некоторые термы в понятии

имеют коэффициенты, значительно превышающие коэффициенты остальных термов — такая закономерность указывает на то, что эти термы часто встречаются вместе в одних и тех же документах и, вероятно, имеют семантическую связь. Пространство, охватываемое векторами понятий, часто называют *латентно-семантическим пространством*.

Строки матрицы \mathbf{V} соответствуют документам, а столбцы — понятиям. Следовательно, любую строку можно интерпретировать как вектор коэффициентов, определяющих значимость соответствующих понятий в данном документе. Эта матрица является двойником исходной матрицы термов/документ в том смысле, что каждый элемент v_{ij} можно рассматривать как *частоту* понятия в документе, так же как каждый элемент матрицы термов/документ отражает частоту терма.

Представление SVD позволяет вычислять сходство между запросами и документами на основе понятий. Во-первых, используя строки матрицы \mathbf{V} , можно вычислить косинусное сходство между документами. А так как каждая строка является векторным представлением соответствующего документа в пространстве понятий, сходство двух документов можно вычислить непосредственно:

$$\cos(\mathbf{v}_i, \mathbf{v}_j) = \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|}. \quad (4.45)$$

Далее мы должны преобразовать запрос в вектор понятий, чтобы вычислить косинусное сходство между запросом и документами. Этот процесс называется *сверткой запросов*. Уравнение 4.44 можно преобразовать для выражения векторов документов в виде функций от матрицы термов/документ:

$$\mathbf{V} = \mathbf{X}^T \mathbf{U} \Sigma^{-1}. \quad (4.46)$$

Пользовательский запрос можно рассматривать как еще один документ, соответствующий некоторому вектору \mathbf{q} частот термов, поэтому его можно подставить в уравнение 4.46 как вырожденный случай матрицы термов/документ с одним столбцом:

$$\mathbf{p} = \mathbf{q}^T \Sigma^{-1}, \quad (4.47)$$

где \mathbf{p} — требуемое представление запроса на основе понятий. Получив это представление, можно оценить соответствие документов запросу, используя косинусное сходство на основе понятий:

$$\text{score}(\mathbf{q}, d_i) = \cos(\mathbf{p}, \mathbf{v}_i) = \frac{\mathbf{p} \cdot \mathbf{v}_i}{\|\mathbf{p}\| \|\mathbf{v}_i\|}. \quad (4.48)$$

Уравнение 4.48 определяет новый метод оценки — оценки латентно-семантического индексирования (Latent Semantic Indexing, LSI), который можно использовать в качестве альтернативы стандартной модели векторного пространства и подходу TF×IDF. Основным преимуществом оценки LSI перед методами на основе термов является возможность извлечения документов, не содержащих термов из запроса. Например, вектор понятий может включать три ключевых термина — *candy* (конфеты), *sweet* (леденцы) и *confection* (сладости), — имеющие сильные семантические зависимости и часто используемые вместе. Документ, содержащий только слова *candy* и *sweet*, все равно будет иметь большой коэффициент для этого понятия в своем векторном представлении. То же верно для запроса *confection*. Следовательно, документ и запрос будут иметь большое косинусное сходство благодаря этому понятию, даже в отсутствие общих термов.

Следующий шаг — уменьшение размерности, то есть ограничение числа размерностей понятий $k < r$. Уменьшение размерности увеличивает ошибку реконструкции 4.43, но в целом этот шаг полезен применительно к латентно-семантическому анализу (LSA), потому что уменьшает шум и оставляет только размерности с самой высокой энергией. Напомним, что SVD гарантирует упорядочение столбцов в матрице \mathbf{U} по их важности¹. То есть понятие \mathbf{u}_1 соответствует наиболее стойкой и частой комбинации термов, тогда как понятие \mathbf{u}_r — наименее значимой комбинации. Таким образом, выполняя свертывание, мы сохраняем только самые сильные понятия и оставляем самые левые столбцы матриц \mathbf{U} и \mathbf{V} , тем самым уменьшая размерность базиса понятий и пространства документов. Важным параметром метода LSA является количество сохраняемых понятий. Часто он выбирается из нескольких возможных значений путем эмпирической оценки точности и полноты. Оптимальное количество понятий, как правило, намного меньше количества разных термов в коллекции; 300–500 понятий — вполне хорошее значение даже для больших коллекций [Bradford, 2008]. Для числа понятий k разложение 4.44 преобразуется следующим образом:

$$\begin{aligned} \mathbf{X}_k &= \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T \\ &= \left[\begin{array}{ccc} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_k \\ | & & | \end{array} \right]_{m \times k} \left[\begin{array}{ccc} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_k \end{array} \right]_{k \times k} \left[\begin{array}{ccc} - & \mathbf{v}_1 & - \\ & \vdots & \\ - & \mathbf{v}_n & - \end{array} \right]^T_{n \times k}. \end{aligned} \quad (4.49)$$

¹ Подробное объяснение значения слова «важность» в контексте SVD вы найдете в главе 2.

Эта операция усечения восстанавливает не точную исходную матрицу термов/документ, а только ее приближение \mathbf{X}_k . Документы по-прежнему соответствуют строкам матрицы \mathbf{V} , но каждый вектор имеет только k элементов. Другими словами, документы и запросы отображаются в пространство с k измерениями, и это пространство используется для вычисления метрики сходства.

ПРИМЕР 4.4

Латентно-семантический анализ довольно трудно понять без численного примера, поэтому в оставшейся части этого раздела мы рассмотрим такой пример. Он довольно мал, но составлен так, чтобы подчеркнуть основные особенности LSA. Однако имейте в виду, что LSA — это метод машинного обучения, требующий значительных объемов данных для эффективного применения на практике. Начнем с коллекции из трех небольших документов о кондитерских магазинах:

d_1 : Chicago Chocolate. Retro candies made with love.

d_2 : Chocolate sweets and candies. Collection with mini love hearts.

d_3 : Retro sweets from Chicago for chocolate lovers.

После фильтрации стоп-слов и применения простой нормализации и стемминга получаем следующую матрицу термов/документ:

$$\mathbf{X} = \begin{matrix} & d_1 & d_2 & d_3 \\ \text{chicago} & \begin{bmatrix} 1 & 0 & 1 \end{bmatrix} \\ \text{chocolate} & \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \\ \text{retro} & \begin{bmatrix} 1 & 0 & 1 \end{bmatrix} \\ \text{candy} & \begin{bmatrix} 1 & 1 & 0 \end{bmatrix} \\ \text{made} & \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \\ \text{love} & \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \\ \text{sweet} & \begin{bmatrix} 0 & 1 & 1 \end{bmatrix} \\ \text{collection} & \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \\ \text{mini} & \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \\ \text{heart} & \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \end{matrix}. \quad (4.50)$$

Применив сингулярное разложение и уменьшив размерность до двух понятий, получим следующие факторные матрицы:

$$\begin{array}{c}
 \text{concept 1} \quad \text{concept 2} \\
 \begin{array}{l}
 \text{chicago} \\
 \text{chocolate} \\
 \text{retro} \\
 \text{candy} \\
 \text{made} \\
 \text{love} \\
 \text{sweet} \\
 \text{collection} \\
 \text{mini} \\
 \text{heart}
 \end{array}
 \begin{bmatrix}
 -0,318 & \mathbf{0,424} \\
 -\mathbf{0,486} & 0,018 \\
 -0,318 & \mathbf{0,424} \\
 -0,333 & -0,148 \\
 -0,166 & 0,257 \\
 -\mathbf{0,488} & 0,018 \\
 -0,320 & -0,239 \\
 -0,168 & -0,406 \\
 -0,168 & -0,406 \\
 -0,168 & -0,406
 \end{bmatrix}
 \end{array}
 \quad (4.51)$$

$$\Sigma_2 = \begin{bmatrix} 3,562 & 0 \\ 0 & 1,966 \end{bmatrix}, \quad (4.52)$$

$$\begin{array}{c}
 \text{concept 1} \quad \text{concept 2} \\
 \begin{array}{l}
 d_1 \\
 d_2 \\
 d_3
 \end{array}
 \begin{bmatrix}
 -0,591 & 0,505 \\
 -0,598 & -0,798 \\
 -0,541 & 0,329
 \end{bmatrix}
 \end{array}
 \quad (4.53)$$

Первое наблюдение, которое можно сделать: столбцы матрицы \mathbf{U}_2 подчеркивают некоторые логические темы, которые можно найти в тексте. Два самых больших коэффициента в первой колонке соответствуют термам *chocolate* (шоколад) и *love* (любовь), за которыми следуют коэффициенты для термов *sweet* (леденцы) и *candy* (конфеты). Наибольшие коэффициенты во втором столбце соответствуют термам *Chicago* (Чикаго) и *retro* (ретро). Это объясняется наличием двух документов, последовательно использующих один и тот же набор слов, говоря о теме *retro* и *Chicago*, и все три документа последовательно используют одни и те же слова, говоря о любви и шоколаде.

Второе наблюдение следует из матрицы документов \mathbf{V}_2 . Первый столбец соответствует понятию Chocolate&Love. Все коэффициенты в столбце имеют одинаковый знак, то есть векторы всех трех документов указывают в одном направлении вдоль этой оси. Второй столбец соответствует понятию Retro&Chicago, но векторы документов указывают в разных направлениях, потому что эта тема упоминается только в первом и третьем документах.

Теперь попробуем запросить документы, используя два запроса: *Chicago* и *candy*. Запросам соответствуют следующие векторы с частотами термов (порядок термов совпадает с матрицей 4.50):

$$\begin{aligned} \mathbf{q}_{\text{chicago}} &= [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0], \\ \mathbf{q}_{\text{candy}} &= [0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0]. \end{aligned} \quad (4.54)$$

Преобразуя эти векторы с помощью формулы 4.47 и вычисляя косинусные сходства с векторами документов из матрицы \mathbf{V}_2 , получаем оценки, представленные в табл. 4.4. Как видите, только первый и третий документы имеют высокие оценки соответствия запросу *Chicago*, что вполне ожидаемо. Второй запрос, *candy*, представляет более интересный случай. Все три документа получают высокую оценку, хотя в третьем документе терм *candy* отсутствует. Это объясняется тем, что *candy* является частью понятия *Chocolate&Love*, явно присутствующем в третьем документе. Метод LSA сумел распознать связь между запросом и документом через это понятие и присвоить документу соответствующий ранг.

Таблица 4.4. Окончательные оценки документов в примере латентно-семантического анализа

Запрос	\mathbf{d}_1	\mathbf{d}_2	\mathbf{d}_3
Chicago	0,891	-0,510	0,806
Candy	0,183	0,969	0,338

Латентно-семантический анализ разрабатывался как альтернатива методам поиска на основе модели векторного пространства, таким как стандартная оценка $\text{TF} \times \text{IDF}$. Эмпирическое исследование показывает, что во многих ситуациях он способен превзойти базовую модель векторного пространства. Кроме того, LSA предлагает следующие преимущества:

СИНОНИМЫ. Представление с меньшим числом измерений способно учитывать синонимы и семантические отношения. LSA также можно использовать для оценки расстояний между словами при создании тезауруса, который затем можно применять для замены синонимов в стандартном методе $\text{TF} \times \text{IDF}$. Кроме того, существуют специализированные методы на основе LSA для вычисления семантического сходства, такие как метод оценки коррелированных вхождений с аналогичной лексической семантикой (Correlated Occurrence Analogue to Lexical Semantic, COALS) [Rohde et al., 2006].

УМЕНЬШЕНИЕ ШУМА. Уменьшение размерности может эффективно устранять шум и избыточность данных.

ВЫСОКАЯ ПОЛНОТА. Поиск на основе LSA успешно работает с запросами и документами, не имеющими общих термов. Это позволяет достичь высокой полноты.

АВТОМАТИЗАЦИЯ. Латентно-семантический анализ опирается на разложение матриц без учителя, соответственно данный анализ легко автоматизируется.

С другой стороны, LSA имеет ряд недостатков, обусловленных, главным образом, его эвристическим характером, который пренебрегает сложными статистическими свойствами текстов:

ПОЛИСЕМИЯ. Латентно-семантический анализ имеет ограниченную способность решать проблему полисемии. Даже при том что LSA может связывать одно и то же слово с несколькими понятиями, фиксируя тот факт, что слово может иметь разные значения в зависимости от контекста, он не в состоянии различать разные значения слова в документе, потому что все значения усредняются до частоты термина в матрице терм/документ. Это ограничение обусловлено природой модели «мешок слов» и не позволяет LSA распознавать более тонкие семантические отношения между словами.

ПОЛНОТА. Теоретическая основа LSA является неполной, потому что не предлагает никакой модели документов и термов.

ИНТЕРПРЕТИРУЕМОСТЬ. Понятия, создаваемые LSA, порой трудно интерпретировать из-за отрицательных значений и отсутствия формальной модели документа.

ДОПУЩЕНИЕ О НОРМАЛЬНОСТИ. Одним из ключевых преимуществ метода главных компонент, используемого в LSA, является возможность создания некоррелированных векторов понятий. Принцип некоррелированности основан на предположении, что данные имеют нормальное (гауссово) распределение, для которого нулевая корреляция между компонентами подразумевает независимость. Однако это допущение неверно для матриц счета, таких как матрица терм/документ.

В следующих нескольких разделах мы попытаемся преодолеть некоторые ограничения LSA. Но сначала обсудим, как можно заменить эвристическую модель разложения надежным вероятностным решением, а затем переиначим подход «мешок слов», чтобы усовершенствовать учет семантических отношений между словами.

4.5.4. Вероятностное тематическое моделирование

Вероятностное тематическое моделирование — это семейство методов семантического анализа, учитывающих семантические отношения между документами и словами с помощью скрытых переменных, называемых *темами* (topic). Одно из главных допущений в тематическом моделировании заключается в том, что документы генерируются терм за термом посредством некоторого вероятностного процесса. Этот процесс моделирует тематическую структуру коллекции документов с помощью скрытых переменных, которые можно интерпретировать как темы. Каждый документ обычно представлен набором тем, и каждая тема определяет распределение слов в документе. Генеративный процесс предназначен только для отражения определенных статистических свойств документов, таких как тематическое распределение слов, и не обязательно создает текст, читаемый или похожий на настоящий. На рис. 4.17 показаны основные компоненты тематической модели.

Затем модель генеративного процесса обучается путем вывода значений латентных переменных из данных. Данные — это набор документов, в котором каждый документ является простой группой термов, а скрытые параметры распределения тем и отношения между темами и термами являются всего лишь абстракциями, которые никогда не наблюдаются непосредственно, но могут быть оценены. Процесс статистического вывода можно рассматривать как подъем снизу вверх на рис. 4.17. Обученная модель описывает связи между термами и темами (какие слова наиболее характерны для данной темы) и отношения между темами и документами (о чем говорит документ). Поиск документов также может осуществляться путем обучения запроса в предполагаемой структуре темы и вычисления сходства между запросом и документами в скрытом пространстве тем.

Метод тематического моделирования напоминает LSA в том смысле, что тоже использует понятие скрытых тем и отображает документы в векторные представления в тематическом пространстве. В то же время его математическая основа сильно отличается от LSA. Это очень важно, потому что данная математическая основа позволяет использовать не одну модель, а целое семейство мощных методов и приемов. В контексте алгоритмического маркетинга эта группа методов важна не только для поиска, но и для служб рекомендаций, потому что обеспечивает общую основу для моделирования отношений между разными сущностями, такими как слова и документы или пользователи и продукты. В следующих разделах мы обсудим две популярные тематические модели — вероятностный латентно-семантический анализ и латентное размещение Дирихле.

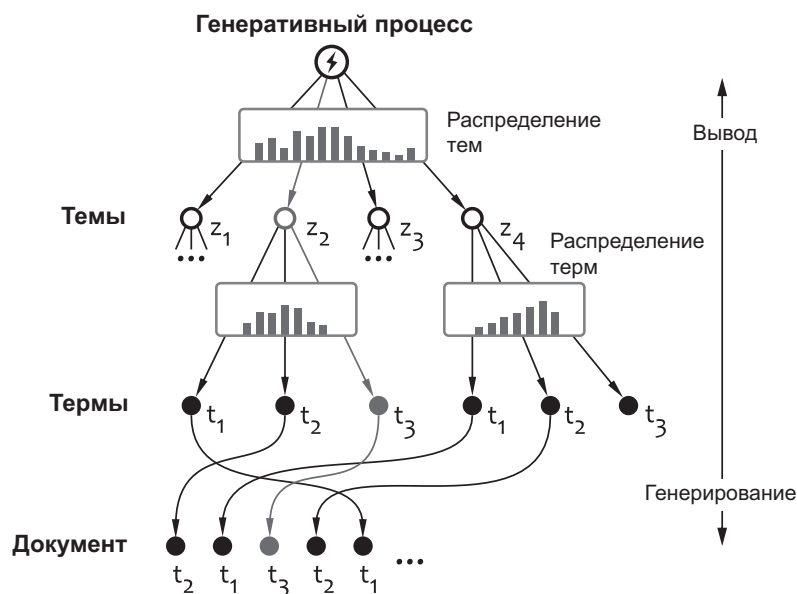


Рис. 4.17. Концептуальное представление вероятностных тематических моделей. Термы документа генерируются последовательно, и каждый терм соответствует некоторому пути в графе модели

4.5.5. Вероятностный латентно-семантический анализ

Вероятностный латентно-семантический анализ (Probabilistic Latent Semantic Analysis, pLSA) — один из основных методов вероятностного тематического моделирования [Hofmann, 1999]. Несмотря на то что он подходит к задаче семантического анализа с вероятностной точки зрения, полученную структуру модели можно рассматривать как матричное разложение, что позволяет напрямую сравнивать pLSA с латентно-семантическим анализом на основе SVD. pLSA можно рассматривать с двух разных точек зрения. Первая — это модель скрытых переменных, то есть вероятностная модель, использующая скрытые переменные (темы) для объяснения отношений между документами и термами. Вторая — разложение матриц, связывающее вероятностную модель скрытых переменных с LSA. Мы обсудим эти два аспекта отдельно в следующих разделах [Oneata, 1999].

4.5.5.1. Модель скрытых переменных

Модель pLSA принадлежит семейству вероятностных тематических моделей. Чтобы описать модель pLSA более формально, определим сначала следующие три основные сущности:

ДОКУМЕНТЫ $D = \{d_1, \dots, d_n\}$ — множество n документов.

ТЕРМЫ $T = \{t_1, \dots, t_m\}$ — множество m термов (слов), содержащее все различные термы из всех документов.

ТЕМЫ $Z = \{z_1, \dots, z_k\}$ — множество k тем, и k — это параметр модели. Понятие «тема» соответствует понятию «скрытая переменная» в LSA.

Мы явно наблюдаем пары документов и термов $\{d_j, t_{ij}\}$, но не темы. Модель скрытых признаков предполагает, что каждый документ может соответствовать нескольким темам, а вероятности термов в документе определяются темой. Например, представим две темы, которые можно найти в каталоге продуктового магазина: *dairy* (молочные продукты) и *desserts* (десерты). Некоторые описания продуктов в каталоге будут относиться в основном к молочным продуктам, некоторые будут связаны с десертами, а некоторые — соответствовать обоим темам, смешанным в определенной пропорции. Хотя темы не наблюдаются непосредственно, распределение термов в документе, соответствующем молочным продуктам, будет определяться соответствующей темой. Эту идею можно выразить более формально, предположив, что документы создаются следующим генеративным процессом:

1. Создать документ d_j из распределения вероятностей $\Pr(d)$.
2. Для каждого терма t_i в документе d_j :
 - 2.1. Выбрать тему z_l из распределения $\Pr(z_l | d_j)$.
 - 2.2. Выбрать терм t_i из распределения $\Pr(t_i | z_l)$.

Этот процесс соответствует вероятностной модели, изображенной на рис. 4.18. Каждый документ моделируется как комплекс тем, а распределение лексем в документе определяется темами. Та же модель, но в более компактном виде, изображена на рис. 4.19.

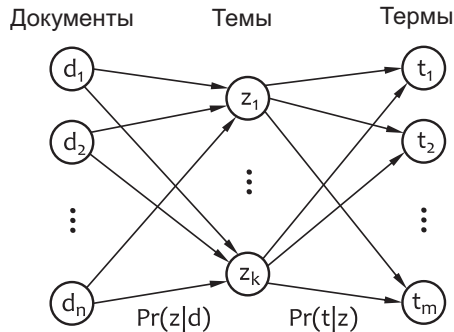


Рис. 4.18. Подробная структура модели pLSA

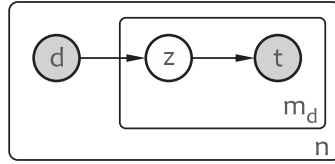


Рис. 4.19. Графическое представление модели pLSA. Внешний блок представляет цикл выбора документов. Внутренний блок представляет цикл выбора тем и термов в документе, содержащем m_d термов. Затененные круги соответствуют наблюдаемым переменным, незатененный — скрытой переменной

По аналогии с LSA, модель pLSA рассматривает каждый документ как мешок слов. С вероятностной точки зрения это означает, что пары документ/терм (d, t) условно независимы:

$$\Pr(D, T) = \prod_{d,t} \Pr(d, t). \quad (4.55)$$

Кроме того, модель pLSA предполагает, что терм и документы условно независимы для данной темы, то есть

$$\Pr(t | d, z) = \Pr(t | z). \quad (4.56)$$

Модель сопряженной вероятности для $D \times T$ можно выразить как

$$\Pr(d, t) = \Pr(d) \Pr(t | d), \quad (4.57)$$

для которого условную вероятность термина в документе можно выразить как сумму вероятностей для всех тем:

$$\begin{aligned} \Pr(t | d) &= \sum_z \Pr(t, z | d) = \sum_z \Pr(t | d, z) \Pr(z | d) = \\ &= \sum_z \Pr(t | z) \Pr(z | d). \end{aligned} \quad (4.58)$$

Подставив выражение 4.58 в 4.57, получим законченное определение модели:

$$\begin{aligned} \Pr(d | t) &= \sum_z \Pr(d) \Pr(t | z) \Pr(z | d) = \sum_z \Pr(d | z) \Pr(t | z) = \\ &= \sum_z \Pr(z) \Pr(t | z) \Pr(d | z). \end{aligned} \quad (4.59)$$

Следующий шаг — изучение ненаблюдаемых вероятностей и вывод скрытых тем. Для набора обучающих документов D функция правдоподобия определяется как

$$L = \Pr(D, T) = \prod_{d,t} \Pr(d, t)^{n(d,t)}, \quad (4.60)$$

где $n(d, t)$ — число появлений термина t в документе d , то есть частота термина. Упрощая функцию правдоподобия взятием логарифма, получаем следующее уравнение:

$$\begin{aligned} \log L &= \sum_{d,t} n(d, t) \cdot \log \Pr(d, t) = \\ &= \sum_{d,t} n(d, t) \cdot \log \sum_z \Pr(z) \Pr(t | z) \Pr(d | z). \end{aligned} \quad (4.61)$$

Вероятности термов $\Pr(t | z)$, вероятности документов $\Pr(d | z)$ и вероятности тем $\Pr(z)$ являются параметрами модели, которая должна быть обучена так, чтобы максимизировать правдоподобие. Это эквивалентно решению следующей задачи оптимизации:

$$\begin{aligned} \log L &\text{ получит максимальное значение} \\ \text{при условии } \sum_t \Pr(t | z) &= 1 \\ \sum_d \Pr(d | z) &= 1 \\ \sum_z \Pr(z) &= 1. \end{aligned} \quad (4.62)$$

Эту задачу можно решить с помощью алгоритма максимизации ожидания (Expectation–Maximization, EM), стандартного подхода к оценке максимального правдоподобия в моделях со скрытыми переменными [Hofmann, 1999]. Проблема, однако, в том, что у нас есть $k(m - 1)$ параметров вероятности $\Pr(t | z)$ для всех возможных пар термов и $k(n - 1)$ параметров вероятности $\Pr(d | z)$ для всех пар документов и тем. Обратите внимание, что число параметров равно $k(m - 1)$, а не km , из-за ограничений нормализации вероятности, описанных задачей 4.62. То есть число параметров велико и растет линейно с размером коллекции документов. Эта проблема считается одним из основных недостатков модели pLSA.

Учитывая оценки параметров, отношения между документами, терминами и темами, а также семантический смысл тем, можно проанализировать путем исследования величин условных вероятностей. Эти же параметры можно проиндексировать и сохранить в качестве поисковых запросов [Park and Ramamohanarao, 2009]. По аналогии с латентно-семантическим анализом сходство между запросом и документом можно вычислить в скрытом семантическом пространстве как косинусное расстояние, или скалярное произведение, между двумя векторными представлениями. В случае pLSA векторное представление запроса q и документа d в латентно-семантическом пространстве задаются условными вероятностями $\Pr(q | z)$ и $\Pr(d | z)$ соответственно. Тогда меру сходства можно определить как следующее скалярное произведение:

$$\text{score}(q, d) = \sum_z \Pr(q|z) \cdot \Pr(d|z). \quad (4.63)$$

Значения $\Pr(d|z)$ известны из модели, но представление запроса $\Pr(q|z)$ необходимо определять для каждого запроса. Этого можно достичь фиксацией параметров $\Pr(t|z)$ и $\Pr(z)$ и обучением модели 4.62 с учетом $\Pr(q|z)$. Метрику сходства затем можно использовать для оценки и ранжирования документов в результатах поиска.

4.5.5.2. Разложение матрицы

Даже притом что подход со скрытыми переменными сильно отличается от LSA (вместо алгебраического разложения матрицы используется вероятностный процесс), эти два метода тесно связаны. Это можно продемонстрировать, переписав модель скрытых переменных с использованием матричной нотации. Во-первых, напомним, что LSA аппроксимирует матрицу частот термов, определяемую выражением 4.41 как произведение трех матриц:

$$\mathbf{X} = \mathbf{U} \cdot \mathbf{\Sigma} \cdot \mathbf{V}^T, \quad (4.64)$$

где \mathbf{X} — матрица $m \times n$ частот термов для всех пар термов и документов, \mathbf{U} — матрица $m \times k$ координат термов в пространстве понятий, и \mathbf{V} — матрица $n \times k$ координат документов в пространстве понятий. С другой стороны, мы определили, что модель сопряженной вероятности в pLSA определяется произведением трех факторов:

$$\Pr(d, t) = \sum_z \Pr(z) \Pr(t|z) \Pr(d|z). \quad (4.65)$$

Если переписать это выражение в матричной нотации, мы получим модель pLSA в форме, которую можно непосредственно сравнить с разложением LSA:

$$\mathbf{P} = \mathbf{L} \cdot \mathbf{S} \cdot \mathbf{R}^T, \quad (4.66)$$

где \mathbf{L} — матрица $m \times k$ с вероятностями всех термов $\Pr(t|z)$, \mathbf{R} — матрица $n \times k$ с вероятностями всех документов $\Pr(d|z)$, и \mathbf{S} — диагональная матрица $k \times k$ априорных вероятностей тем $\Pr(z)$. Другими словами, pLSA, подобно LSA, можно рассматривать как алгоритм разложения матриц, но в данном случае разложение управляется другой целью. Если в LSA целью является минимизация ошибки аппроксимации, то в pLSA целью является максимизация функции правдоподобия или, как вариант, минимизация расхождения между наблюдаемым распределением и моделью.

4.5.5.3. Свойства pLSA

Модель pLSA предлагает несколько важных преимуществ перед LSA. Во-первых, направления в пространстве pLSA неотрицательны и интерпретируются как веро-

ятности. Направления в пространстве LSA не имеют формальной интерпретации, а значения, полученные разложением LSA, могут быть отрицательными, что также усложняет интерпретацию.

Вторым важным отличием является подход к решению проблемы *полисемии*. LSA способен связывать синонимы, находящиеся рядом в скрытом семантическом пространстве, но обычно не различает разные значения одного и того же слова в зависимости от контекста. pLSA, напротив, распределяет вероятностную массу термина по нескольким темам, которые могут соответствовать разным значениям слова [Hofmann, 1999]. В частности, если один и тот же терм t наблюдается в двух разных документах, d_i и d_j , тема, с которой он имеет самую сильную ассоциацию в контексте первого документа:

$$\operatorname{argmax}_z \Pr(z | d_i, t),$$

может отличаться от темы, с которой этот же терм ассоциируется в контексте второго документа

$$\operatorname{argmax}_z \Pr(z | d_j, t).$$

Несмотря на эти преимущества, pLSA обычно имеет более сложную реализацию, чем LSA. Если LSA основан на детерминированном разложении SVD, то для оценки параметров модели в pLSA требуется итерационный алгоритм максимизации ожиданий. Модель pLSA также имеет несколько структурных проблем, которые мы обсудим и рассмотрим в следующем разделе.

4.5.6. Латентное размещение Дирихле

Модель pLSA является важным шагом вперед по сравнению с LSA. Она устанавливает прочную статистическую основу, которая позволяет расширять, упрощать и объединять различные модели с использованием вероятностных методов. Однако модель pLSA имеет несколько недостатков:

- Каждый документ представлен вектором вероятностей, а не генеративной вероятностной моделью. Эти вероятности являются параметрами, которые должны определяться по имеющимся данным. В результате получается большое количество параметров, которое линейно растет с количеством термов и документов и увеличивает риск переобучения.
- pLSA не накладывает ограничений на связи документов и термов с темами. Интуитивно мы ожидаем, что каждый документ будет связан с небольшим количеством тем и каждая тема будет связана с небольшим количеством тер-

мов, но pLSA не предлагает явных параметров для управления этим аспектом модели.

Эти проблемы можно решить, создав модель с более сложным генеративным процессом, чем процесс pLSA, описанный выше. В этом разделе мы рассмотрим один из самых ярких примеров таких моделей — латентное размещение Дирихле (Latent Dirichlet Allocation, LDA). Модель LDA можно рассматривать как обобщение pLSA. Она является одной из самых популярных и широко используемых вероятностных тематических моделей [Blei et al., 2003]. Модель LDA основана на понятии распределения Дирихле, о котором подробнее рассказывается в приложении А в конце книги.

Так же как pLSA, модель LDA использует подход на основе скрытых переменных, который предполагает, что каждый документ соответствует комплексу скрытых тем, а термины документа подчиняются распределениям, связанным с темами [Blei et al., 2003]. Если допустить, что число скрытых тем k предопределено, модель LDA можно описать с помощью следующего генеративного процесса для каждого документа d из коллекции документов D :

1. Получить некоторое количество термов в документе m_d из некоторого случайного распределения. Выбор распределения не критичен для архитектуры модели.
2. Получить k -мерный вектор вероятностей \mathbf{p} из распределения Дирихле $Dir(\alpha)$, где α является параметром модели. Каждый элемент вектора \mathbf{p} интерпретируется как вероятность соответствующей темы, то есть этот вектор определяет смесь тем.
3. Для каждого термина в документе:
 - 3.1. Выбрать тему z_i согласно вектору вероятностей \mathbf{p} , то есть $\Pr(z_i = i | \mathbf{p}) = p_i$.
 - 3.2. Выбрать терм t из полиномиального распределения вероятностей $\Pr(t | z_i; \beta)$, зависящий от темы z_i . Это распределение определяется как параметр модели β для каждой пары термина и темы.

Основное отличие от процесса pLSA, описанного в разделе 4.5.5.1, заключается в том, что модель LDA получает темы из глобального параметрического распределения, а не из распределений, полученных из документов. Параметрами этой модели являются k -мерный параметр Дирихле α и $k \times m$ матрица вероятностей термов β , где m — общее число различных термов во всех документах. Строки матрицы β определяют полиномиальные распределения слов для соответствующих тем. Эти параметры выбираются один раз для коллекции документов, благодаря чему число параметров получается меньше, чем в pLSA. На рис. 4.20 изображено графическое представление генеративного процесса.

В контексте одного документа сопряженное распределение комплекса тем, всех тем и всех термов задается как

$$\Pr(\mathbf{p}, \mathbf{z}, \mathbf{t}) = \Pr(\mathbf{p} | \boldsymbol{\alpha}) \prod_t \Pr(t | z_t; \beta) \Pr(z_t | \mathbf{p}), \quad (4.67)$$

где распределение $\Pr(\mathbf{p}; \boldsymbol{\alpha})$ определяется как $Dir(\boldsymbol{\alpha})$ и параметры $\boldsymbol{\alpha}$ и β заданы. Обратите внимание, что $\Pr(z_t | \mathbf{p})$ — это просто значение вероятности из \mathbf{p} , соответствующее z_t . Частное распределение документа можно получить интегрированием по вероятностям тем и суммированием по всем темам:

$$\Pr(d) = \int \Pr(\mathbf{p} | \boldsymbol{\alpha}) \prod_t \sum_z \Pr(t | z_t; \beta) \Pr(z_t | \mathbf{p}) d\mathbf{p}. \quad (4.68)$$

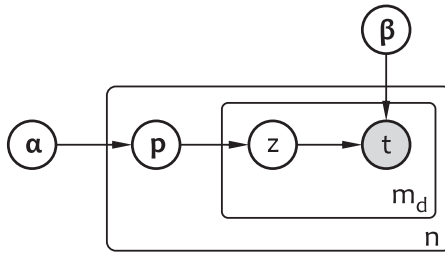


Рис. 4.20. Графическое представление модели LDA

Функцию правдоподобия для коллекции документов можно получить путем вычисления произведения вероятностей документа:

$$\Pr(D) = \prod_d \int \Pr(\mathbf{p}_d | \boldsymbol{\alpha}) \prod_{t \in d} \sum_z \Pr(t | z_d; \beta) \Pr(z_d | \mathbf{p}_d) d\mathbf{p}_d. \quad (4.69)$$

Процесс обучения этой модели усложняется тем, что параметры $\boldsymbol{\alpha}$ и β связаны внутренней суммой в уравнении 4.69. Эту проблему можно решить с помощью методов вывода апостериорного распределения, таких как вариационный вывод и семплирование по Гиббсу [Blei et al., 2003; Asuncion et al., 2009].

Модель LDA решает две проблемы pLSA, о которых упоминалось выше. Во-первых, она уменьшает число параметров, определяя другой генеративный процесс, не использующий параметры, зависящие от документов. Во-вторых, предшествующее априорное распределение Дирихле формирует вероятности тем, накладывая штраф за отношения между темами и документами.

4.5.7. Модель Word2Vec

Word2Vec — это семейство моделей, созданное с целью преодолеть ограничения методов семантического анализа, основанных на представлении «мешок слов», путем

учета локального контекста слова, а не всего документа [Mikolov et al., 2013a, b]. Двумя основными типами моделей Word2Vec являются модели непрерывного мешка слов и словосочетаний с пропуском (skip-gram). Суть подхода на основе *непрерывного мешка слов* заключается в построении предиктивной модели, оценивающей вероятность слова на основе одного или нескольких слов в окружающем контексте, как показано на рис. 4.21. Слова в скользящем окне контекста интерпретируются как мешок слов, то есть учитываются только отдельные термины и их частоты, а не их порядок. Предиктивную модель можно построить так, что каждое слово будет ассоциироваться с вектором весов, определяемых в процессе обучения модели. Эти векторы затем можно интерпретировать как представления слов в некотором скрытом семантическом пространстве, подобно векторам, создаваемым в латентно-семантическом анализе или тематическими моделями, то есть такое векторное представление можно использовать для поиска и создания тезауруса. Подход на основе словосочетаний с пропуском (skip-gram) является противоположностью модели непрерывного мешка слов — он принимает целевое слово на входе и предсказывает контекст. Однако архитектуры предиктивных моделей непрерывного мешка слов и словосочетаний с пропуском очень похожи, поэтому в оставшейся части раздела все свое внимание мы сосредоточим на первой из них.

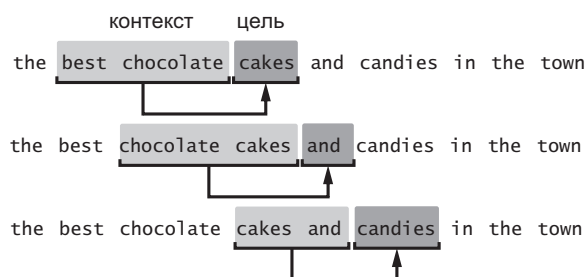


Рис. 4.21. Пример модели непрерывного мешка слов с окном контекста, включающим два слова

Для выявления семантических отношений и прогнозирования термина на основе контекста модель Word2Vec использует неглубокую нейронную сеть. Сначала мы обсудим архитектуру сети в предположении, что в контексте есть только одно слово, а затем обобщим результат на случай с несколькими словами. Нейронная сеть, используемая в модели Word2Vec, состоит из входного, скрытого и выходного слоев, изображенных на рис. 4.22.

На вход сети подается бинарный вектор, представляющий контекст. Если общее количество термов в коллекции равно n , то входной вектор имеет n элементов и каждый элемент равен единице, если соответствующий терм присутствует в кон-

тексте, и ноль в противном случае. Поскольку мы рассматриваем случай контекста из одного слова, обозначим единственный контекстный терм как t_k . Тогда входной вектор будет иметь только один ненулевой элемент x_k :

$$x_i = \begin{cases} 1, & i = k \\ 0, & \text{иначе} \end{cases}. \quad (4.70)$$

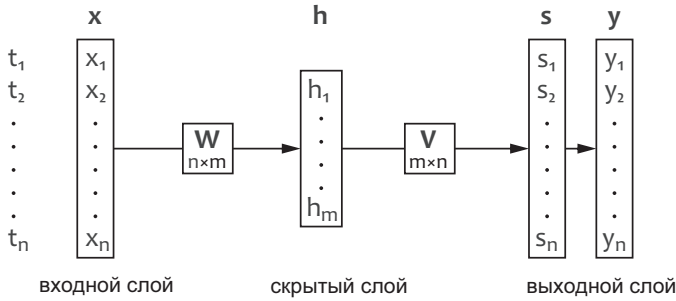


Рис. 4.22. Архитектура нейронной сети в модели Word2Vec для контекста с одним словом

Входной вектор преобразуется в m промежуточных выходов с использованием узлов скрытого слоя. Это преобразование выбрано линейным (каждый промежуточный выход h_i является взвешенной суммой входов x_j), то есть определяется через матрицу весов W :

$$W_{n \times m} = \begin{bmatrix} \text{---} & w_1 & \text{---} \\ & \vdots & \\ \text{---} & w_n & \text{---} \end{bmatrix}. \quad (4.71)$$

Теперь промежуточные выходы можно выразить как произведение матрицы весов на входной вектор. Согласно условию, входной вектор содержит только один ненулевой элемент, поэтому результат будет идентичен соответствующей строке в матрице весов:

$$h = W^T x = w_k^T. \quad (4.72)$$

Оценка на выходе создается комбинацией линейного преобразования и функции softmax. Линейная часть, по аналогии со скрытым слоем, определяется с помощью весовой матрицы V :

$$V_{m \times n} = \begin{bmatrix} | & & | \\ v_1 & \cdots & v_n \\ | & & | \end{bmatrix}. \quad (4.73)$$

Эта матрица умножается на промежуточные выходы, чтобы получить оценку для каждого из n членов:

$$s_i = \mathbf{v}_i^T \mathbf{h}, \quad i = 1, \dots, n. \quad (4.74)$$

Сигналы s_i являются произвольными значениями, но мы должны интерпретировать их как предсказанные вероятности соответствующих термов в данном контексте. Другими словами, мы решаем задачу многоклассовой классификации, где контекст требуется отнести к одному из n классов, соответствующих предсказанному терму. Как говорилось в главе 2, стандартным способом отображения вектора произвольных значений в вероятности категорий является функция softmax, поэтому конечные результаты определяются следующим образом:

$$y_i = \Pr(t_i | t_k) = \frac{\exp(s_i)}{\sum_{j=1}^n \exp(s_j)}. \quad (4.75)$$

Сеть, определенную выше, можно обучить с использованием стандартных методов обучения искусственных нейронных сетей. Мы не будем углубляться в детали алгоритмов обучения, но кратко рассмотрим основные шаги, чтобы получить представление, как обучаются модели Word2Vec [Rong, 2014b]. Модель обучается итеративно, принимая образцы пар контекста и целевого слова, пропуская контекст через сеть, сравнивая вывод сети с целью и корректируя весовые коэффициенты в матрицах \mathbf{W} и \mathbf{V} . Предположим, что в данной итерации фактически наблюдаемым целевым термом для контекста t_k является t_a . Согласно принципу максимального правдоподобия, наша цель состоит в том, чтобы максимизировать предсказанную вероятность фактического терма с учетом контекста (потому что в конечном итоге нам нужно максимизировать математическое ожидание этой вероятности во всех контекстах):

$$\max \mathbb{E}_{t_k, t_a} [\Pr(t_a | t_k)]. \quad (4.76)$$

Максимизация этой вероятности эквивалентна минимизации следующей функции потерь:

$$J = -\log \Pr(t_a | t_k). \quad (4.77)$$

Подставляя определение 4.75 выхода сети в определение 4.77 функции потерь, находим:

$$J = -\log y_a = -s_a + \log \sum_{j=1}^n \exp(s_j). \quad (4.78)$$

Наша цель — минимизировать функцию потерь по весам w и v . Это можно сделать, изменяя веса методом стохастического градиентного спуска на основе ошибок прогнозирования. Стратегия заключается в том, чтобы начать с выходной стороны сети и рассчитывать изменения весов для матрицы \mathbf{V} , исходя из наблюдаемых ошибок прогнозирования. Затем переместиться на слой назад и вычислить изменения весов для матрицы \mathbf{W} . Этот подход известен как обратное распространение ошибок, или просто *обратное распространение*. В каждом слое мы должны вычислить градиент функции потерь относительно весов. Делается это в два этапа — сначала вычисляется градиент по отношению к оценкам, а затем результат используется для вычисления градиента относительно весов. Итак, начнем с производной по оценкам выходного слоя:

$$\begin{aligned}\frac{\partial J}{\partial s_j} &= -\mathbb{I}(j=k) + \frac{\partial}{\partial s_j} \log \sum_{i=1}^n \exp(s_i) = \\ &= -\mathbb{I}(j=k) + \frac{\exp(s_j)}{\sum_{i=1}^n \exp(s_i)} = \\ &= y_j - \mathbb{I}(j=k) = \\ &= e_j,\end{aligned}\tag{4.79}$$

где $\mathbb{I}(j=k)$ — индикаторная функция, равная единице, если $j=k$, и нулю в противном случае. Как видите, эта производная является просто ошибкой предсказания, поэтому обозначим ее как e_j . Взяв производную по весам выходного слоя, найдем градиент для оптимизации веса:

$$\frac{\partial J}{\partial v_{ij}} = \frac{\partial J}{\partial s_j} \cdot \frac{\partial s_j}{\partial v_{ij}} = e_j \cdot h_i.\tag{4.80}$$

Этот результат означает, что мы должны уменьшить вес v_{ij} , если произведение $e_j \cdot h_i$ положительно, и увеличить в противном случае. Уравнение стохастического градиентного спуска для весов будет выглядеть следующим образом:

$$\mathbf{v}_j^{(new)} = \mathbf{v}_j^{(old)} - \lambda \cdot e_j \cdot \mathbf{h}, \quad j=1, \dots, n,\tag{4.81}$$

где λ — параметр скорости обучения. Следующий шаг — повторить процесс для скрытого слоя. Сначала возьмем производную функции потерь по промежуточным выходам:

$$\frac{\partial J}{\partial h_i} = \sum_{j=1}^n \frac{\partial J}{\partial s_j} \cdot \frac{\partial s_j}{\partial h_i} = \sum_{j=1}^n e_j \cdot v_{ij} = \varepsilon_i.\tag{4.82}$$

Результат, обозначенный как ε , можно интерпретировать как взвешенную сумму ошибок прогнозирования. Это значение нужно вычислить для каждого из m скрытых узлов, то есть получить m -мерный вектор ошибок прогнозирования:

$$\varepsilon = [\varepsilon_1, \dots, \varepsilon_m]. \quad (4.83)$$

Далее вычислим градиент относительно весов скрытого слоя:

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_{ji}} = \varepsilon_i \cdot x_j. \quad (4.84)$$

Мы используем этот результат и стохастический градиентный спуск для изменения весов скрытого слоя, аналогично методу изменения весов выходного слоя для уравнений 4.80 и 4.81. Принимая во внимание тот факт, что все значения x_j в уравнении 4.84 являются нулями, кроме x_k , мы должны обновить только k -ю строку матрицы \mathbf{W} :

$$\mathbf{w}_k^{(new)} = \mathbf{w}_k^{(old)} - \lambda \cdot \varepsilon^T. \quad (4.85)$$

Модель Word2Vec можно обучить применением уравнений 4.81 и итеративно 4.85 для обучения пар контекстов и целевых слов. Этот процесс, однако, требует большого объема вычислений, потому что в соответствии с уравнением 4.81 требуется обновить векторы весов \mathbf{v} для всех термов в каждой обучающей выборке, а число термов n может быть большим. Это требует применения в практических реализациях модели Word2Vec методов оптимизации, таких как иерархическая функция softmax и негативное семплирование (negative sampling) [Mikolov et al., 2013b; Rong, 2014b].

Полученные результаты легко обобщить для случая контекста с несколькими словами. Входной вектор для контекста из q слов, то есть q ненулевых элементов, можно рассматривать как нормализованную сумму (то есть среднее) q контекстов, каждое с одним словом. Это решение иллюстрирует рис. 4.23, хотя фактическая структура сети не изменяется.

Это позволяет переписать уравнение для скрытого слоя следующим образом:

$$\mathbf{h} = \frac{1}{q} \mathbf{W}^T (\mathbf{x}_1 + \dots + \mathbf{x}_q) = \frac{1}{q} (\mathbf{w}_{k_1} + \dots + \mathbf{w}_{k_q})^T. \quad (4.86)$$

Уравнение функции потерь остается тем же, даже при том что представляет другую условную вероятность:

$$J = -\log \Pr(t_a | t_{k_1}, \dots, t_{k_q}) = -s_a + \log \sum_{j=1}^n \exp(s_j). \quad (4.87)$$

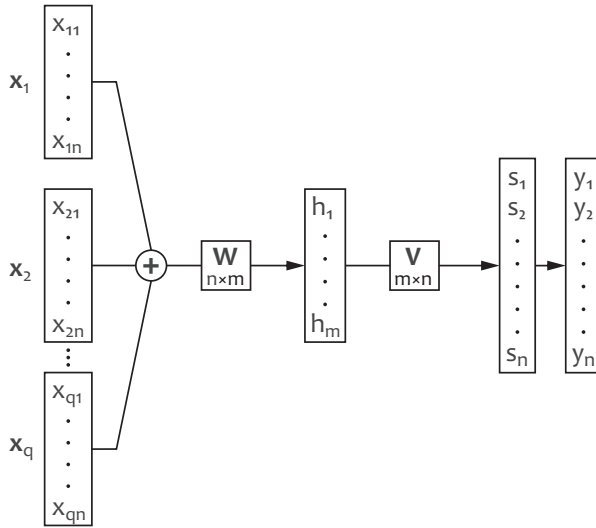


Рис. 4.23. Модель Word2Vec для контекста с несколькими словами

Рассмотрев все вычисления градиента, можно обнаружить, что все уравнения остаются неизменными и дают одну и ту же пару формул обновления весов:

$$\mathbf{v}_j^{(new)} = \mathbf{v}_j^{(old)} - \lambda \cdot e_j \cdot \mathbf{h}, \quad j = 1, \dots, n, \quad (4.88)$$

$$\mathbf{w}_j^{(new)} = \mathbf{w}_j^{(old)} - \frac{\lambda}{q} \cdot e^T, \quad j = 1, \dots, q. \quad (4.89)$$

Единственное отличие — обновляется несколько векторов с весами \mathbf{w} , потому что контекст содержит несколько термов.

ПРИМЕР 4.5

После обучения сети каждому из n термов в коллекции будет соответствовать пара m -мерных векторов, \mathbf{w} и \mathbf{v} . Мощь модели Word2Vec обусловлена тем, что эти векторы обеспечивают представление слов, сохраняющее семантические отношения. Проиллюстрируем это на примере модели Word2Vec, обученной на следующих образцах — парах контекстных и целевых слов [Rong, 2014a]:

drink coffee	tea drink
drink juice	juice drink
drink tea	coffee drink
eat cake	pie coffee

eat pie	cookie juice
eat cookie	cake tea
pie tea	cake coffee

Для захвата семантических шаблонов в этих образцах было решено использовать сеть с 8 скрытыми узлами, поэтому после обучения модели каждый терм был представлен двумя 8-мерными векторами. Для визуализации весовые векторы можно спроецировать в двумерное пространство с помощью метода главных компонент. Например, векторы весов из скрытого слоя w проецируются на плоскость, как показано на рис. 4.24. Как видите, слова сгруппировались в соответствии с шаблонами их использования и, в конечном счете, смыслом.

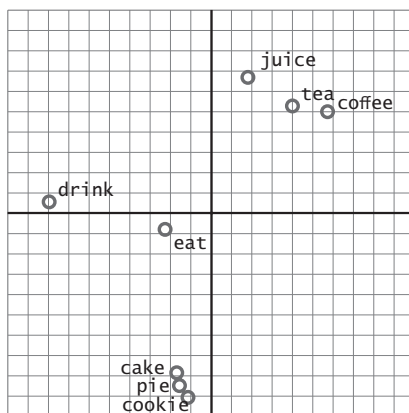


Рис. 4.24. Пример кластеризации слов с помощью модели Word2Vec

Что особенно удивительно, модель Word2Vec, обученная на большой коллекции текстов, создает векторы, поддерживающие своеобразную алгебру в скрытом семантическом пространстве. Рассмотрим следующие примеры, где $v(\cdot)$ обозначает векторное представление терма, полученное моделью Word2Vec, а знак «минус» используется для обозначения стандартной поэлементной разности между векторами:

$$\begin{aligned}
 v(\text{steak}) - v(\text{beef}) &\approx v(\text{salad}) - v(\text{tomato}), \\
 v(\text{steak}) - v(\text{beef}) &\approx v(\text{bread}) - v(\text{flour}), \\
 v(\text{franceq}) - v(\text{paris}) &\approx v(\text{japan}) - v(\text{tokyo}).
 \end{aligned}
 \tag{4.90}$$

Как видите, в первых двух случаях разность векторов отражает идею приготовления пищи и в последнем — идею отношения страна/столица. Другими

словами, прибавив к помидору (tomato) разность между говядиной (beef) и стейком (steak), которую можно интерпретировать как акт приготовления, мы получим салат (salad). Этот тип семантических отношений называется *аналогией слов*. Обратите внимание, что каждое понятие, такое как «приготовление пищи» или «столица», то есть одна из разностей векторов в примерах выше, также является вектором в семантическом пространстве, поэтому можно сделать вывод, что одно направление в пространстве соответствует приготовлению пищи, другое — сжатию страны до ее столицы и т. д.

Одним из применений методов Word2Vec в поиске для продвижения является создание тезауруса. Например, они использовались в системе поиска вакансий Dice.com для преодоления конфликтов в названиях и описаниях вакансий, обусловленных синонимией [Hughes, 2015]. Преимущество этого подхода заключается в том, что он позволяет создать тезаурус, или кластеры слов, пригодный для использования в стандартном механизме синтаксического поиска.

Завершим обзор Word2Vec кратким сравнением с методами LDA, pLSA и LSA. Главное отличие заключается в том, что для выявления семантических зависимостей Word2Vec использует окно локального контекста, тогда как методы тематического моделирования — глобальную статистику документа. Оба подхода имеют свои достоинства и недостатки. Например, в общем случае Word2Vec лучше захватывает аналогии слов, чем приемы тематического моделирования [Pennington et al., 2014]. С другой стороны, векторные представления, создаваемые моделью Word2Vec, не являются разреженными и могут содержать отрицательные элементы, что усложняет их интерпретацию в отличие, например, от результатов латентного размещения Дирихле, способного создавать разреженные векторы, интерпретируемые как вероятности. Это делает Word2Vec менее применимым в приложениях тематического анализа.

4.6. Методы поиска для продвижения

До сих пор мы рассматривали относительно универсальные методы поиска и их применение в поиске для продвижения. Однако задача продвижения товаров в результатах поиска выходит далеко за рамки настройки стандартных методов и требует создания более специализированных методов поиска в таких областях, как электронная коммерция. Это объясняется рядом причин [Khudnev, 2013]:

- *Структурированные сущности*. Многие стандартные методы предназначены для поиска в документах с относительно простой структурой, включающей

одно или несколько текстовых полей. Поиск товаров часто имеет дело с высокоструктурированными документами, более напоминающими записи в реляционных базах данных, чем простой текст. Например, типичный документ с описанием товара может включать сотни числовых и категориальных атрибутов:

Brand: Tommy Hilfiger
Type: Jeans
Color: Black
Weight: Super Skinny
...

Кроме того, часто бывает так, что товар группируется во вложенные сущности или ассоциируется с иерархиями категорий. Например, ретейлер может продавать набор посуды как один товар, но этот товар включает несколько продуктов, и каждый из них, в свою очередь, может иметь несколько вариантов цветов или размеров. Это требует, чтобы служба поиска работала со вложенными или взаимосвязанными сущностями, которые нельзя правильно представить в виде простых документов.

- *Разнообразие товаров.* Для повышения точности результатов многие приложения поиска используют ранжирование. К сожалению, ранжирование имеет ограниченную применимость в поиске для продвижения. Одна из основных причин в том, что результаты, полученные с помощью стандартных методов оценки и смешивания сигналов, как правило, чрезмерно разнообразны, что вызывает отрицательные впечатления у пользователей. Например, запросу *red dress* (красное платье) может соответствовать широкий спектр товаров, содержащих эти два термина в своих атрибутах, включая платья, обувь и даже часы. Продвинутое проектирование и смешивание сигналов, рассмотренные выше, могут помочь улучшить результаты, но они едва ли обеспечат надежное решение проблемы. Другая причина заключается в том, что $TF \times IDF$ и другие популярные методы оценки могут плохо работать со структурированными документами с множеством категориальных полей. И это неудивительно, потому что эти методы рассчитаны на естественные тексты.
- *Составные и полисемичные термины.* Опыт, накопленный в индустрии, показывает, что качество поиска для продвижения существенно зависит от правильной обработки составных и полисемичных терминов. Поисковые запросы в приложениях поиска товаров часто содержат многословные названия брендов и понятия, такие как *Calvin Klein* и *dress shoes*, которые четко передают цель поиска, если обрабатывать их как фразы, но могут неверно истолковываться, если разбить их на отдельные слова. Кроме того, многие названия брендов содержат распространенные слова, что еще больше усложняет правильную интерпретацию. Например, запросу *pink sweater* (розовый свитер) могут со-

ответствовать все товары, произведенные брендом *Pink Rose*, и, наоборот, запросу *pink rose sweater* (свитер от Pink Rose) могут соответствовать все товары розового цвета и цвета розы, хотя он явно указывает на намерение найти определенный бренд (*Pink Rose*).

Наблюдения, перечисленные выше, предполагают необходимость разработки методов поиска, ориентированных на точность, точное соответствие и структуры атрибутов, а не на статистические оценки. Другими словами, мы должны рассмотреть методы поиска, которые интерпретируют документы скорее как записи в базе данных, а не как произвольные тексты, требующие оценки. Ряд таких методов был разработан в компании Масу, ведущим американским ретейлером, для своих служб поиска [Kamotsky and Vargas, 2014; Peter and Eugene, 2015]. Оставшуюся часть раздела мы посвятим обзору этих методов.

4.6.1. Комбинаторный фразовый поиск

Наша первая цель — повысить точность результатов поиска с учетом того, что документы имеют много категориальных полей, которые часто содержат составные и полисемичные термы. Если бы мы могли заставить пользователей писать структурированные логические запросы, это было бы отличным решением данной проблемы. Например, текстовый запрос *pink rose sweater* станет гораздо менее двусмысленным, если пользователь явно сформулирует поля и составные термы:

```
brand:[pink rose] AND type:[sweater]
```

Такой подход можно использовать в некоторых приложениях поиска для продвижения, если дать пользователю удобный интерфейс, помогающий определить отдельные поля. Например, сайты по продаже автомобилей часто предлагают раскрывающиеся списки производителей, моделей автомобилей и других свойств, чтобы пользователь мог задать критерии поиска на уровне полей. Это может оказаться разумным решением для бизнеса с относительно связным ассортиментом, таким как автомобили или недвижимость, но запросы в виде произвольного текста, возможно, предпочтительнее для сфер с большим разнообразием товаров, таких как универмаги.

Запрос с произвольным текстом не содержит ни полей документа, ни границ составных термов, что порождает неоднозначность. Идея комбинаторного фразового поиска состоит в том, чтобы восстановить часть этой информации из исходного запроса с произвольным текстом путем создания нескольких логических запросов с разными комбинациями полей и термов и найти документы, соответствующие этим искусственным запросам. Цель алгоритма, генерирующего запросы, состоит в том, чтобы создать ограниченное число критериев поиска, чтобы найти докумен-

ты, сильно коррелирующие с запросом. Это увеличивает вероятность включения документов в список результатов не из-за случайного совпадения отдельных термов, а благодаря хорошему соответствию атрибутов документа термам и фразам из запроса. Данную методологию можно рассматривать как обобщение методов поиска с n -граммами для документов с несколькими полями.

Первый шаг в комбинаторном фразовом поиске — разбиение запроса на подфразы. Предположим, что пользователь ввел запрос, представляющий собой последовательность из n термов:

$$q = [t_1 \ t_2 \ \dots \ t_n]. \quad (4.91)$$

Существует 2^{n-1} возможных вариантов деления этого запроса на подфразы, потому что между термами в запросе имеется $n - 1$ пробелов, и мы свободны в выборе, делить или не делить запрос по любому пробелу. Например, для запроса с тремя термами есть четыре возможных варианта деления (здесь квадратные скобки используются для обозначения подфраз):

$$\begin{aligned} & [t_1 \ t_2 \ t_3] \\ & [t_1][t_2 \ t_3] \\ & [t_1 \ t_2][t_3] \\ & [t_1][t_2] [t_3]. \end{aligned} \quad (4.92)$$

Второй шаг — создать логический запрос для каждого способа деления так, чтобы каждая подфраза соответствовала одному из полей в документе. Для m подфраз s_1, \dots, s_m в данном варианте деления и документа с k полями f_1, \dots, f_k , логический запрос будет выглядеть следующим образом:

$$\begin{aligned} & (f_1 = s_1 \text{ OR } f_2 = s_1 \text{ OR } \dots \text{ OR } f_k = s_1) \\ & \text{AND } (f_1 = s_2 \text{ OR } f_2 = s_2 \text{ OR } \dots \text{ OR } f_k = s_2) \\ & \dots \\ & \text{AND } (f_1 = s_m \text{ OR } f_2 = s_m \text{ OR } \dots \text{ OR } f_k = s_m). \end{aligned} \quad (4.93)$$

Знак равенства в запросе 4.93 обозначает точное соответствие между подфразовым запросом и значением поля; обе стороны должны точно совпадать, при этом для предварительной обработки запроса и полей можно использовать нормализацию, удаление стоп-слов или стемминг. Запрос 4.93 гарантирует высокую степень охвата данного деления документом в том смысле, что каждая подфраза должна точно совпадать с одним из полей.

Наконец, выполняются логические запросы для всех вариантов деления, и окончательный набор результатов получается как объединение результатов всех логических запросов. Это равносильно объединению всех частичных запросов в один большой логический запрос с помощью оператора OR (ИЛИ). Общая структура запроса представлена на рис. 4.25. Наш алгоритм деления на разделы не пытается распознать составные термины в запросе и механически разбивает их на подфразы. Следовательно, подфразы часто будут нарушать границы составных термов. Например, запрос *blue calvin klein jeans* можно разбить на подфразы *blue calvin* и *klein jeans*. Объединяя логические запросы, мы гарантируем, что хотя бы некоторые из вариантов деления правильно захватят составные термы.

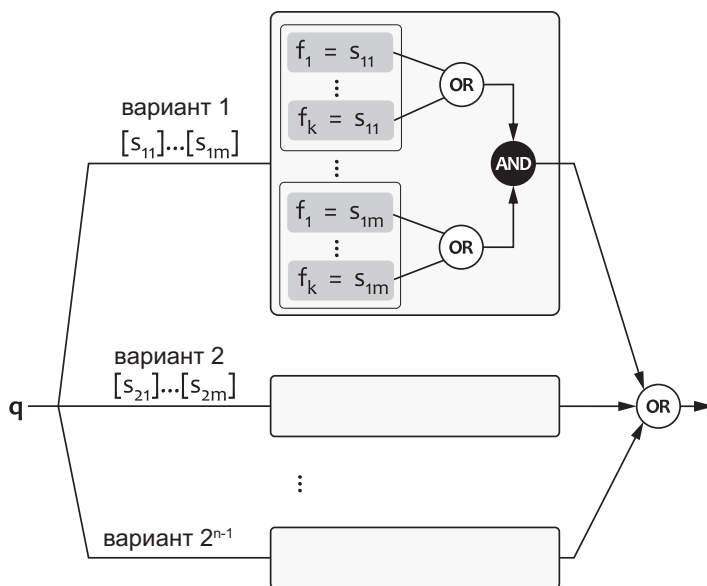


Рис. 4.25. Общая структура запроса в комбинаторном фразовом поиске

ПРИМЕР 4.6

Проиллюстрируем логику работы комбинаторного фразового поиска на примере. Возьмем запрос *pink rose sweater*, который можно разделить четырьмя способами:

- Способ 1 : [pink rose sweater]
- Способ 2 : [pink rose] [sweater]
- Способ 3 : [pinks] [rose sweater]
- Способ 4 : [pinks [rose] [sweater]

Предположим, что товары в каталоге представлены в виде документов с тремя полями: бренд, тип товара и цвет. Запрос в комбинаторном фразовом поиске, сгенерированном для таких полей и вариантов деления, будет иметь структуру, изображенную на рис. 4.26. С увеличением числа полей и термов запрос может получиться очень большим и сложным в вычислительном отношении, но эту проблему можно немного смягчить введением некоторых упрощений. Например, можно ограничить максимальную длину подфраз, поскольку слишком длинные подфразы вряд ли будут представлять значимые составные термы.

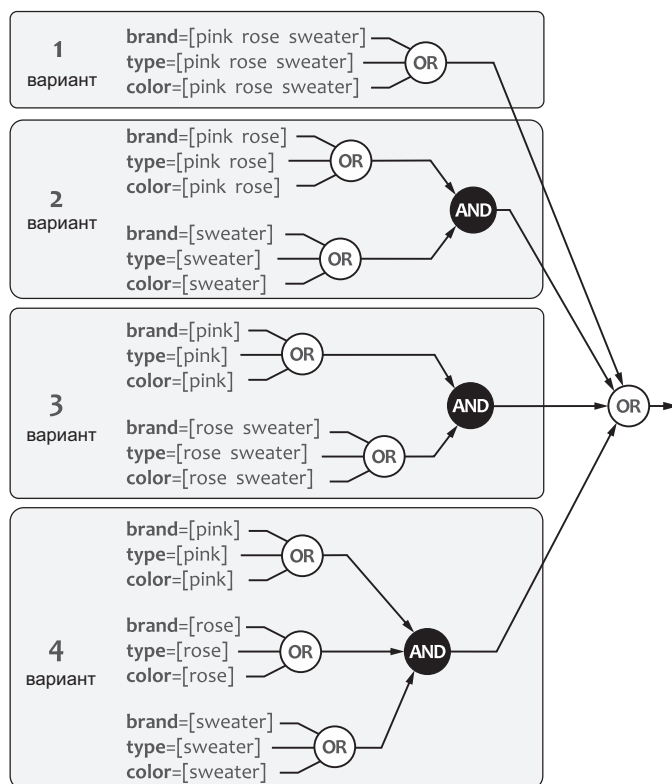


Рис. 4.26. Пример запроса в комбинаторном фразовом поиске

Следующий документ будет соответствовать комбинаторному запросу, потому что поля с названием бренда и типом товара охватывают все подфразы во втором варианте деления:

Brand: pink rose
Type: sweater
Color: black

В то же время товар с типом *sweater* (свитер) и цветом *pink* (розовый) не будет соответствовать запросу, если название бренда будет отличаться от *rose*. Кроме того, комбинаторный фразовый поиск становится еще более ограниченным с ростом длины запроса, потому что должны быть охвачены все термы. Такое поведение отличается от стандартной модели векторного пространства, которая оценивает каждое соответствие термина, из-за чего точность результатов уменьшается с увеличением длины запроса.

Недостаток комбинаторного фразового поиска заключается в экспоненциальном увеличении числа вариантов деления и операторов в окончательном логическом запросе с ростом числа термов в пользовательском запросе. На практике сложность логического запроса часто можно уменьшить, исключив некоторые операторы, опираясь на тип поля. Например, поле *color* (цвет) может иметь ограниченное число допустимых значений, поэтому генератор запросов мог бы отбросить варианты

```
Color = [sweater]  
Color = [pink rose]
```

как бессмысленные.

Комбинаторный фразовый поиск можно рассматривать как метод поиска документов, обеспечивающих полное покрытие запроса с точки зрения подфраз и полей. Однако этот метод можно связать и с семантическим поиском. Хотя комбинаторный поиск не обнаруживает семантических отношений так же явно, как LSA или Word2Vec, он все же пытается выявить и сопоставить составные термы, которые с большой долей вероятности представляют логические понятия. Иначе говоря, его можно рассматривать как попытку сконструировать семантический поиск с помощью примитивов синтаксического поиска [Giunchiglia et al., 2009; Khludnev, 2013].

4.6.2. Контролируемое снижение точности

Комбинаторный фразовый поиск позволяет достичь высокой точности результатов, исключая все документы, которые не полностью покрывают запрос. Это помогает поддерживать согласованность результатов поиска и эффективно использовать доступное экранное пространство. Однако комбинаторный подход имеет свои недостатки. Одна из наиболее серьезных проблем заключается в том, что из-за увеличенной строгости сопоставления такой запрос может возвращать пустой список результатов, особенно если запрос содержит слова с опечатками или иную неудачную комбинацию термов, которые не могут быть охвачены до-

ступными документами. Такое поведение очень нежелательно, потому что вместо списка товаров пользователь увидит пустой экран, а это уменьшит вероятность его конверсии.

Эту проблему можно решить, предприняв дополнительные действия, если основной алгоритм комбинаторного фразового поиска вернет пустой список с результатами. Например, можно сначала попытаться выполнить комбинаторный поиск, требующий точного соответствия полей, а затем вернуться к базовой модели векторного пространства, допускающей частичное соответствие. Эту идею можно развить дальше и создать цепочку методов поиска с постепенно уменьшающейся точностью, вызывать последовательно каждый метод, пока не будет найден хотя бы один документ (или другое минимальное количество документов). Например, такая цепочка может иметь следующую структуру:

1. *Точное соответствие.* Поиск документов с помощью стандартного комбинаторного фразового поиска без нормализации или стемминга.
2. *Нормализация, стемминг и коррекция правописания.* Если список результатов пуст, к полям и термам в запросе применяются нормализации и стемминг и снова запускается комбинаторный фразовый поиск. К термам в запросе также можно применить коррекцию правописания.
3. *Применение n -грамм.* Если совпадений нет, вместо поиска точного совпадения выполните поиск с использованием n -грамм.
4. *Частичное совпадение.* Если совпадений по-прежнему нет, попробуйте повторить поиск, удалив из запроса одно-два слова, чтобы найти документы с частичным покрытием запроса.

Процесс можно прервать на любом этапе и вернуть найденные результаты. Например, для запроса с ошибкой *Abibas sneakers* комбинаторный фразовый едва ли найдет хоть один документ, но для исправленного запроса *Adidas sneakers* почти наверняка найдется достаточно много документов, чтобы остановить дальнейшее ослабление критериев поиска. Этот метод, называемый контролируемым снижением точности, помогает управлять балансом между высокой точностью комбинаторного фразового поиска и риском разочарования пользователя от отсутствия результатов.

4.6.3. Вложенные сущности и динамическая группировка

С точки зрения продвижения товаров, проблему поиска можно рассматривать как проблему эффективного использования экранного пространства. Безусловно, очень важно дать пользователю набор релевантных товаров, отвечающих цели поиска,

но не менее важно продемонстрировать в лучшем виде доступный ассортимент, учитывая, что пользователь готов просмотреть только ограниченное количество результатов. Эффективное использование экранного пространства является одной из главнейших задач в поиске товаров, потому что каталоги часто содержат тесно связанные товары, которые пользователь может воспринимать как дубликаты. Например, иногда разумнее представить пользователю набор релевантных, но разных платьев, а не несколько вариантов одного и того же платья, отличающихся цветом и размером.

Проблема дублирования и неэффективного использования экранного пространства возникает из-за иерархических связей между товарами в каталоге. Характер и структура этих связей сильно зависят от сферы бизнеса. Например, в универсамах часто можно увидеть следующую иерархию продуктов:

- Наименьшей единицей товара является вариант продукта, обычно называемый *единицей складского учета*, или *артикулом*. Примером артикула могут служить джинсы Levi's 511 белого цвета размера 30W × 32L. Все физические экземпляры одного и того же варианта считаются идентичными.
- Продукт — это логическая сущность, включающая один или несколько вариантов продукта. Например, джинсы Levi's 511 — это продукт, включающий варианты разных размеров и цветов. Продукт обычно имеет цену, в том смысле что все его варианты имеют одинаковую цену.
- Несколько продуктов могут быть объединены в *коллекцию продуктов*. Коллекция может быть продана как единое целое, или пользователю может быть представлена возможность выбрать из нее отдельные элементы. Например, несколько тарелок, мисок и кружек могут продаваться как набор посуды по одной цене. При этом пользователь может купить подмножество конкретных вариантов продукта.

Все методы поиска, затрагивавшиеся до сих пор, предполагают, что элементы каталога смоделированы как простые документы, поэтому мы должны отобразить иерархическую структуру с коллекциями продуктов, продуктами и их вариантами в коллекцию документов. Один из возможных способов — представить каждый вариант продукта как отдельный документ, чтобы каждый элемент в списке с результатами поиска соответствовал одному варианту продукта. В общем случае такой подход вполне оправдан и широко применяется на практике, но он страдает проблемой дубликатов и, как следствие, влечет неэффективное использование экранного пространства. Эту проблему иллюстрирует рис. 4.27: прием моделирования документов для вариантов дает формально релевантный результат, но не эффективен с точки зрения продвижения в сравнении с моделированием на уровне продукта, которое лучше демонстрирует имеющийся ассортимент.

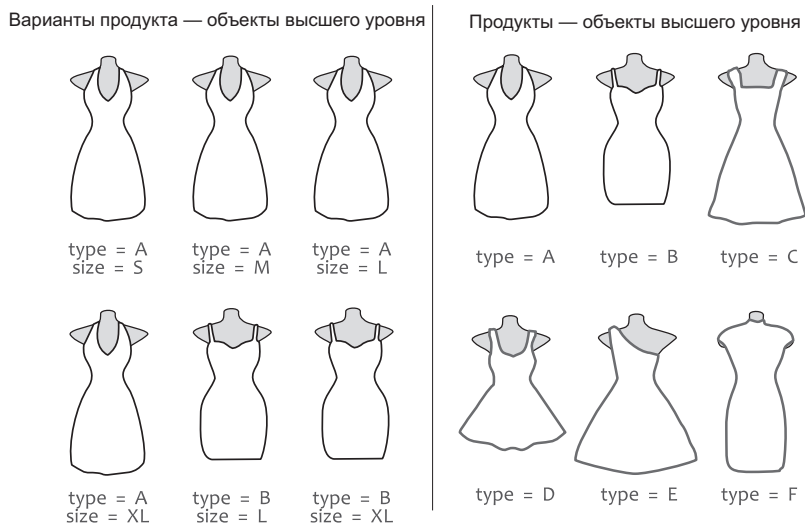


Рис. 4.27. Примеры результатов поиска для запроса *evening dresses* (вечерние платья) с разными подходами к моделированию данных

Моделирование документов на уровне продукта может решить проблему с дубликатами, но влечет свои собственные проблемы. Самый простой подход к моделированию на уровне продукта — представить каждый продукт в виде единого документа. Для этого все атрибуты продукта и варианта продукта нужно объединить в простой список полей; иначе говоря, продукт должен наследовать все атрибуты своих вариантов. Например, рассмотрим два варианта дорожного чемодана, относящихся к одному продукту:

Brand: Samsonite
 Name: Carry-on Hardside Suitcase
 Color: red
 Size: small

Brand: Samsonite
 Name: Carry-on Hardside Suitcase
 Color: black
 Size: large

Два варианта можно объединить в один документ со следующей структурой:

Brand: Samsonite
 Name: Carry-on Hardside Suitcase
 Color: red black
 Size: small large

Результат выглядит разумным, потому что документ получает высокую оценку для таких запросов, как *red suitcase* (красный чемодан), *small suitcase* (маленький чемодан) и т. д. Основная проблема такого подхода в том, что он теряет струк-

турную информацию о вложенных сущностях, из-за чего невозможно отличить допустимые и недопустимые комбинации атрибутов. Документ выше получит хорошую оценку для запроса *small red suitcase* (маленький красный чемодан), и это правильно, потому что один из вариантов действительно является маленьким и красным чемоданом. Но тот же документ получит столь же высокую оценку для запроса *small black suitcase* (маленький черный чемодан). А это уже неверно, поскольку ни один из вариантов не является маленьким и черным одновременно, что делает продукт нерелевантным для данного запроса. Эта проблема довольно сложна с точки зрения реализации, потому что ее нельзя решить исключительно с позиции простых документов, и требуется, чтобы механизм поиска либо явно поддерживал вложенные сущности, либо мог обрабатывать документы уровня вариантов и обобщать результаты, группируя варианты в продукты. Если фильтрация продуктов реализована правильно, результаты уровня продуктов могут существенно повысить эффективность службы поиска товаров.

Мы видели, что объединение вариантов продуктов в продукты может принести определенную пользу, поэтому далее подумаем о возможности объединения продуктов в коллекции. Этот вопрос еще сложнее, потому что пользователь может преследовать разные цели и искать либо продукты, либо коллекции продуктов. Например, пользователь, ищущий по запросу *dinnerware* (набор посуды), вероятнее всего, ожидает получить результаты с коллекциями, а пользователь, ищущий по запросу *cup* (чашка), — список отдельных продуктов. Значит, мы должны динамически принимать решение о группировке, опираясь на запрос и найденные результаты. С этой целью можно ввести эвристические правила для анализа структуры результатов и совпавших атрибутов и принятия решения о группировке. Например, отдельные продукты можно заменить коллекцией, если коллекция в целом соответствует запросу, то есть если все или почти все продукты в коллекции и атрибуты самой коллекции соответствуют запросу. Рассмотрим пример на рис. 4.28. Запросу *white cup* (белая чашка), скорее всего, будут соответствовать отдельные продукты или коллекции, включающие только белые чашки, но не наборы посуды с тарелками (plate), мисками (bowl) или чашками (cup) разных цветов. Следовательно, мы должны вернуть пользователю результаты, содержащие в основном отдельные продукты. С другой стороны, поиск по запросу *white dinnerware* (набор белой посуды), скорее всего, имеет целью получить другой результат. Можно ожидать результатов со значительным количеством наборов посуды, состоящих в основном из белых предметов, соответствующих терму *white* (белый) на уровне вариантов и терму *dinnerware* (набор посуды) на уровне коллекций. Такие наборы посуды хорошо соответствуют запросу, поэтому их можно включить в список результатов как коллекции, а не как отдельные продукты.

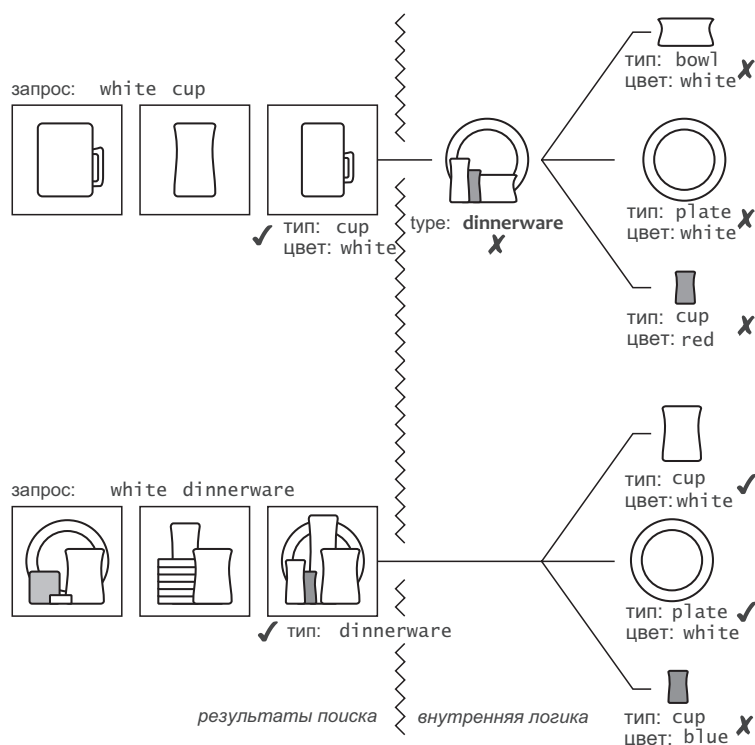


Рис. 4.28. Пример динамической группировки

4.7. Настройка релевантности

Мы рассмотрели множество способов сопоставления, ранжирования и смешивания сигналов, обеспечивающих широкие возможности управления релевантностью и продвижением товаров. Тонкая настройка всех этих средств управления для согласованной работы может оказаться очень сложной задачей, требующей разработки методов и приемов оптимизации. Мы уже сделали шаг в этом направлении, введя несколько метрик качества поиска, таких как точность, полнота и взвешенная накопленная релевантность. Следующим нашим шагом станет разработка методов настройки релевантности с использованием этих метрик в роли целей оптимизации.

Напомню, что стандартные метрики релевантности можно использовать для оценки качества одного списка результатов или среднего качества результатов по набору запросов. Целью оптимизации является максимизация общих экономических показателей службы поиска, то есть мы должны оптимизировать метрики релевантности для набора запросов с наибольшим вкладом в общую выгоду от

использования службы. Иначе говоря, общую экономическую эффективность службы поиска можно неформально выразить как сумму вкладов отдельных поисковых запросов:

$$\text{Revenue} = \sum_{q \in Q} R(q) \cdot m(q), \quad (4.94)$$

где q — запрос, Q — набор возможных запросов, $R(q)$ — средний доход от конверсий, обусловленных запросом q , и $m(q)$ — метрика качества результатов поиска, которая, как предполагается, должна быть пропорциональна количеству конверсий. На практике едва ли можно оптимизировать все возможные запросы, но мы можем отобрать из исторических данных наиболее популярные запросы, приносящие доход, и использовать для оптимизации только этот набор. Процесс оптимизации релевантности можно организовать как непрерывную оценку и улучшение средней эффективности по этому набору запросов. Он состоит из следующих шагов:

1. Собирается и анализируется статистика использования службы для выявления набора запросов с наибольшим вкладом в эффективность службы. Будем называть эти запросы эталонными.
2. Для эталонных запросов вычисляются метрики релевантности, чтобы оценить общую эффективность службы поиска.
3. Вручную анализируются результаты поиска для каждого эталонного запроса и производится настройка алгоритмов поиска для улучшения метрик релевантности.
4. Новые настройки сначала тестируются на подмножестве реальных пользователей, а затем применяются постоянно.

Описанный процесс можно повторять снова и снова, получая отзывы пользователей об изменениях в алгоритмах и поддерживая набор эталонных запросов в актуальном состоянии. Расчет метрик релевантности и настройку алгоритмов поиска можно рассматривать как узкое место в программном конвейере, поскольку оба этапа требуют участия человека для принятия решений о релевантности и доработке формул оценки. Можно попробовать устранить этот пробел, разработав методы автоматической настройки формул и оценивать релевантность результатов поиска, анализируя поведение пользователей и их взаимодействие со службой поиска. Следующие разделы мы посвятим изучению этих двух тем.

4.7.1. Обучение ранжированию

Основная цель службы поиска состоит в ранжировании документов в соответствии с их релевантностью заданному запросу и контексту. Интуитивно понятно, что эта

задача тесно связана с задачами классификации или регрессии — имея определенный запрос, нужно точно предсказать *степень релевантности*, или *ранг*, документа, а затем сконструировать список результатов, отсортировав документы в соответствии с предиктивными оценками. Эта задача, которую часто называют обучением ранжированию, активно исследуется учеными в области информационного поиска и компаниями, специализирующимися на веб-поиске, такими как Yahoo и Microsoft. Это привело к появлению большого числа исследовательских работ, отчетов о промышленном использовании и наборов тестовых данных для оценки и сравнения методов обучения ранжированию. Хотя число алгоритмов обучения ранжированию довольно велико, многие из них используют схожие методы проектирования признаков и целевые функции, которые можно рассматривать как общую основу для обучения ранжированию. Кроме того, некоторые конкурсы, такие как Learning to Rank Challenge, организованный компанией Yahoo в 2010 году, показали, что сложные методы обучения ранжированию имеют весьма ограниченное преимущество на реальных данных перед более простыми [Chapelle and Chang, 2011]. Учитывая это, мы сосредоточимся на обобщенной основе обучения ранжированию и в качестве примера рассмотрим один конкретный алгоритм. Желаящие могут найти обширный каталог алгоритмов обучения ранжированию в Liu, 2009.

Задачу обучения ранжированию формально можно определить следующим образом. Существует обучающий набор, содержащий Q образцов, где каждый образец представлен парой, включающей поисковый запрос и соответствующий список результатов. Список результатов для запроса q содержит документы m_q , и каждому документу d в списке присваивается степень релевантности $y_{q,d}$. Предположим, что y — категориальная переменная, принимающая одно из K значений. Например, набор оценок релевантности может включать пять значений: 1 — *идеально*, 2 — *отлично*, 3 — *хорошо*, 4 — *удовлетворительно* и 5 — *плохо*. Если определить функцию, преобразующую пару с запросом q и документом d в вектор признаков $x_{q,d}$, обучающий набор можно естественным образом представить в виде коллекции векторов признаков и соответствующих обучающих меток:

$$\begin{aligned} (x_{q,d}, y_{q,d}), \quad q = 1, \dots, Q \\ d = 1, \dots, m_q. \end{aligned} \tag{4.95}$$

На практике обучающий набор данных можно создать, извлекая списки результатов для каждого запроса с помощью обычных методов поиска и устанавливая оценки релевантности, исходя из суждений экспертов. Цель состоит в том, чтобы обучить модель ранжирования, которая предсказывает оценку y по входным данным, состоящим из запроса и документа.

Как и другие приемы обучения с учителем, обучение ранжированию начинается с проектирования признаков. Как уже упоминалось, оценка релевантности про-

гнозируется для документа в контексте определенного запроса, поэтому вектор признаков зависит как от документа, так и от запроса. На практике обычно используются следующие группы признаков [Chapelle and Chang, 2011; Liu and Qin, 2010].

ПРИЗНАКИ ДОКУМЕНТА. Признаки этого типа содержат статистики и атрибуты документа, включая:

- Основные статистические характеристики документа, такие как количество термов. Эти данные можно рассчитывать независимо для каждого поля и для всего документа и таким образом получить несколько групп признаков.
- Классифицирующие метки, такие как тип продукта, ценовая категория и т. д.
- Динамические атрибуты и веб-статистика. Примерами таких признаков могут служить данные о продажах, рейтинги пользователей и новизна.
- Реализация веб-поиска обучения ранжированию часто включает веб-графы и признаки, описывающие аудиторию, такие как количество входящих и исходящих ссылок для веб-страницы. Даже имея ограниченную применимость в поиске товаров, такие показатели при их доступности вполне можно считать неплохими кандидатами на роль признаков.

ПРИЗНАКИ ЗАПРОСА. К признакам этого вида относятся различные статистики, связанные с запросом. Как и признаки документа, эту группу можно разбить на несколько подкатегорий:

- Простые статистики запросов, например количество термов.
- Статистики использования запросов, такие как частота использования и процент переходов по ссылке.
- Атрибуты, полученные из набора результатов, связанного с запросом. Например, запрос можно связать с темой, такой как *мебель*, если большая часть результатов относится к этой категории.

ПРИЗНАКИ ДОКУМЕНТА/ЗАПРОСА. Признаки, связанные как с запросом, так и с документом. Это наиболее важная категория признаков и может включать следующие группы:

- Различные статистики, рассчитанные для термов, присутствующих и в запросе, и в документе. Например, это может быть сумма или дисперсия частот термов или обратных частот документов для общих термов. Эти метрики можно рассчитать для каждого поля документа, а также для всего документа.
- Стандартные метрики соответствия и сходства текста, такие как количество общих термов и $TF \times IDF$.

- Статистики, связанные с отзывами пользователей. Сюда можно отнести вероятности разных взаимодействий, такие как вероятность щелчка (доля пользователей, щелкнувших на данном документе хотя бы один раз среди всех пользователей, которые ввели определенный запрос), вероятность последнего щелчка (доля пользователей, завершивших поиск на данном документе), вероятность пропуска (доля пользователей, щелкнувших на документе ниже данного) и т. д.

Структура вектора признаков изображена на рис. 4.29. Общее количество признаков в практических приложениях может достигать нескольких сотен.

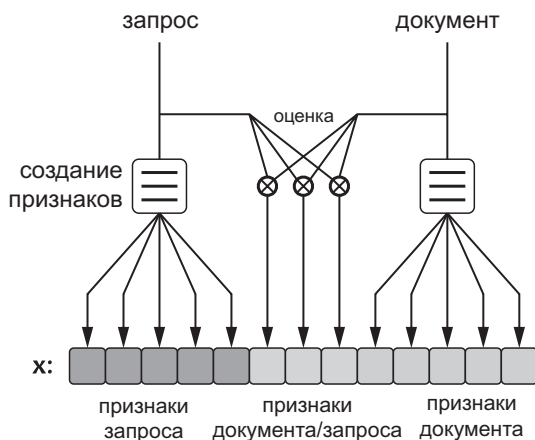


Рис. 4.29. Проектирование признаков для обучения ранжированию

Следующий шаг в создании модели ранжирования — определение функции потерь, которая будет использоваться при обучении модели в качестве цели. Хотя обучение ранжированию тесно связано с классификацией и регрессией, определение функции потерь — не самая простая задача, потому что для нас важно расположить документы в результатах в относительном порядке, не совпадающем со стандартными ошибками классификации или регрессии. Один из возможных подходов — определить функцию потерь как меру релевантности, такую как взвешенная накопленная релевантность (Discounted Cumulative Gain, DCG). К сожалению, DCG является невыпуклой и негладкой функцией, что может быть проблемой для многих алгоритмов обучения с учителем, основанных на градиентном подходе. Хотя DCG часто используется для оценки качества обучения алгоритмов ранжирования, большинство методов используют для обучения разные функции потерь. Эти функции обычно делятся на три категории:

ПОТОЧЕЧНЫЕ. Поточечный подход пытается предсказать степень релевантности каждого документа независимо и тем самым свести ранжирование к стандартной

задаче регрессии или классификации. Следовательно, общая функция потерь L_0 определяется как сумма ошибок прогнозирования для отдельных степеней:

$$L_0 = \sum_{q,d} L(f(\mathbf{x}_{q,d}), y_{q,d}), \quad (4.96)$$

где $f(\mathbf{x}_{q,d})$ — прогнозируемая степень релевантности, а $L(\cdot)$ — функция потерь классификации или регрессии. Потери классификации, например, могут определяться с помощью индикаторной функции, равной нулю, если прогноз правильный, и единице в противном случае:

$$L(f(\mathbf{x}_{q,d}), y_{q,d}) = \mathbb{I}(f(\mathbf{x}_{q,d}) \neq y_{q,d}). \quad (4.97)$$

Однако функция потерь классификации также является невыпуклой и негладкой, поэтому может потребоваться аппроксимировать ее с помощью некоторой другой функции. Мы обсудим возможные варианты далее в этом разделе вместе с конкретными алгоритмами обучения ранжированию. Другой альтернативой является использование функции потерь регрессии:

$$L(f(\mathbf{x}_{q,d}), y_{q,d}) = (y_{q,d} - f(\mathbf{x}_{q,d}))^2. \quad (4.98)$$

Можно показать, что ошибка DCG ограничена сверху потерями классификации и регрессии, поэтому минимизация функций потерь помогает оптимизировать DCG [Cossock and Zhang, 2006; Li et al. 2007,]. Однако поточечный подход имеет существенный недостаток, независимо от выбора функции потерь. Дело в том, что нас интересует относительный порядок результатов в списке, а не качественные или количественные оценки отдельных степеней. Например, поточечный подход не признает, что мы получим идеально ранжированный список результатов из четырех элементов с оценками релевантности (1, 2, 3, 4), даже если оценки предсказаны как (2, 3, 4, 5). Следовательно, мы можем по-другому взглянуть на функцию потерь, чтобы учесть относительный порядок элементов.

Поточечный подход используется во многих алгоритмах ранжирования, включая McRank [Li et al., 2007] и PRank [Crammer and Singer, 2001].

ПОПАРНЫЕ. Попарный подход пытается преодолеть ограничения поточечных методов, накладывая штраф не за неправильно предсказанные оценки релевантности, а за пары документов, ранжированные в обратном порядке. То есть общая функция потерь определяется как сумма функций попарных потерь для всех пар документов в списке результатов с разными оценками:

$$L_0 = \sum_q \sum_{i,j: y_{q,i} > y_{q,j}}^{m_q} L(f(\mathbf{x}_{q,i}), f(\mathbf{x}_{q,j})). \quad (4.99)$$

Функция попарных потерь часто определяется на основе разности между прогнозируемыми оценками, соответственно документы, ранжированные в обратном порядке, вносят свой вклад в потери. Например, функцию можно определить как экспоненциальные потери:

$$L(f(\mathbf{x}_{q,i}), f(\mathbf{x}_{q,j})) = \exp(f(\mathbf{x}_{q,j}) - f(\mathbf{x}_{q,i})). \quad (4.100)$$

Попарный подход также можно рассматривать как задачу классификации, но, в отличие от поточечной классификации, он направлен на классификацию пар документов (пар, где первый документ более релевантен, чем второй, и пар, где более релевантен второй документ).

Примерами попарных алгоритмов ранжирования могут служить RankNet [Burges et al., 2005], RankBoost [Freund et al., 2003] и RankSVM [Herbrich et al., 2000].

ПОСПИСОЧНЫЙ. В посписочном подходе функция потерь определяется на основе всего списка результатов. Иначе говоря, роль «экземпляров» для обучения в посписочном методе играют списки документов, а не отдельные документы или их пары, как в поточечном и попарном подходах. Функция потерь имеет довольно общую форму, которая принимает список пар с предиктивными и фактическими оценками релевантности:

$$L_0 = \sum_q L\left(\left(f(\mathbf{x}_{q,1}), y_{q,1}\right), \dots, \left(f(\mathbf{x}_{q,m_q}), y_{q,m_q}\right)\right). \quad (4.101)$$

Внутреннюю функцию потерь L можно определить как меру релевантности, такую как DCG, или ее гладкую аппроксимацию. К методам посписочного ранжирования относятся AdaRank [Xu and Li, 2007] и ListRank [Cao et al., 2007].

Итак, мы рассмотрели подготовку обучающих данных и описали возможные варианты функции потерь. Осталось сделать последний шаг — выбрать модель прогнозирования и обучить ее, минимизируя потери прогноза. Поскольку обучение ранжированию тесно связано с классификацией, для ранжирования можно адаптировать многие стандартные методы обучения с учителем. В частности, опыт Yahoo и Microsoft показал, что на практике особенно эффективны деревья решений, нейронные сети и их ансамбли [Chapelle and Chang, 2011; Burges, 2010]. Мы завершаем этот раздел обзором алгоритма McRank, который для прогнозирования оценок релевантности использует увеличивающиеся деревья решений [Li et al., 2007].

McRank — это поточечный алгоритм обучения ранжированию, который сводит задачу ранжирования к множественной классификации. Как говорилось выше, оценки релевантности являются категориальными переменными с K классами:

$$y_{q,d} \in \{1, 2, \dots, K\}. \quad (4.102)$$

Наша цель — создать модель классификации, оценивающую вероятность каждого класса на основе вектора признаков:

$$p_{q,d,k} = \Pr(y_{q,d} = k \mid \mathbf{x}_{q,d}), \quad k = 1, \dots, K. \quad (4.103)$$

После определения вероятностей алгоритм McRank сортирует документы в соответствии с их *ожидаемой релевантностью*:

$$r_{q,d} = \sum_{k=1}^K k \cdot p_{q,d,k}. \quad (4.104)$$

Модель классификации создается в McRank с помощью алгоритма дерева градиентного бустинга. Так как это градиентный метод, необходима гладкая функция потерь. McRank использует следующую сглаженную версию ошибки классификации из выражения 4.97:

$$\sum_{q,d} \sum_{k=1}^K -\log(p_{q,d,k}) \mathbb{I}(y_{q,d} = k). \quad (4.105)$$

McRank использует стандартный алгоритм дерева градиентного бустинга, который итеративно создает ансамбль деревьев решений для минимизации функции потерь 4.105. Результатом является модель, оценивающая вероятности, описанные в уравнении 4.103, которую можно использовать для ранжирования документов в списке результатов поиска.

4.7.2. Обучение ранжированию на неявной обратной связи

Обучение ранжированию дает мощную возможность автоматической настройки релевантности, что помогает исключить или упростить смешивание сигналов вручную. Это особенно важно для программных систем. Обучение ранжированию, однако, зависит от мнений экспертов, на основе которых определяются оценки релевантности, используемые в обучении модели. Этот шаг часто требует значительных человеческих усилий, а также ограничивает способность системы автоматически настраиваться динамически. Можно попробовать обойти эту проблему, автоматически определяя оценки релевантности на основе взаимодействий пользователей с результатами поиска. Например, результаты, остающиеся без внимания, скорее всего, являются нерелевантными. Один из возможных способов использования этой информации — ее включение в векторы признаков, как мы

уже делали это в предыдущем разделе. Можно сделать еще шаг вперед и попробовать разработать метод, который определяет оценки релевантности по неявной обратной связи.

Совершенно понятно, что пользователи, как правило, щелкают на релевантных результатах и пропускают нерелевантные, однако поведение пользователя может передавать более сложные отношения релевантности. Например, пользователь может ввести запрос, просмотреть результаты, щелкнуть на некоторых из них, переформулировать запрос и щелкнуть на некоторых новых результатах. Все запросы и документы в этом случае связаны с одной целью поиска, поэтому отношения релевантности можно определять как в пределах одного списка результатов, так и между запросами. В этом разделе мы рассмотрим модель обратной связи, фиксирующую такие отношения с помощью нескольких эвристических правил [Radlinski and Joachims, 2005]. Эта конкретная модель уходит корнями в академические исследования, однако о похожих методах обучения с неявной обратной связью сообщала компания Yahoo [Zhang et al., 2015].

Модель, которую мы рассмотрим, имеет две группы правил интерпретации обратной связи. Первая, изображенная на рис. 4.30, включает два правила, которые распространяются на рамки одного поискового запроса. Первое правило гласит, что если пользователь щелкает на каком-либо документе в списке результатов, этот документ более релевантен для данного запроса, чем все вышестоящие. Оно основано на предположении, что обычно пользователь читает результаты сверху вниз. Второе правило основано на эмпирических наблюдениях (включая результаты наблюдений за движением глаз), что пользователь обычно просматривает не менее двух лучших результатов в списке, прежде чем предпринять какое-либо действие. То есть если пользователь щелкает на первом документе в списке, он считается более релевантным, чем второй (в отношении данного запроса).

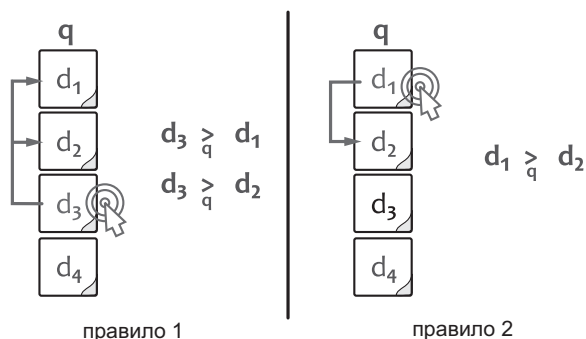


Рис. 4.30. Правила интерпретации неявной обратной связи для документов в одном списке результатов

Вторая группа правил применяется к *цепочкам запросов*, то есть последовательностям запросов с разными формулировками одной и той же цели поиска. Для этого прежде необходимо определить, принадлежат ли запросы одной цепочке. Это нетривиальная задача, потому что пользователь может сделать несколько запросов, преследуя одну цель, но сформулировать их по-разному или сделать несколько несвязанных друг с другом запросов, пытаясь найти совершенно разные продукты. Рассматриваемая нами модель неявной обратной связи подходит к этой проблеме через построение дополнительного классификатора, предсказывающего принадлежность пары запросов одной цепочке. Модель обучается на основе пар запросов, классифицированных вручную, и использует такие функции, как интервал времени между запросами, количество общих термов и количество общих документов в списках результатов. После группировки запросов в цепочки можно ввести четыре дополнительных правила релевантности для применения к парам списков с результатами поиска. Все эти правила основаны на предположении, что запросы в цепочке выражают одну и ту же цель и, следовательно, могут считаться эквивалентными.

Первые два правила в этой группе изображены на рис. 4.31. Они повторяют правила для одного запроса, рассмотренные выше, но применяются к смежным запросам в цепочке. Рассмотрим цепочку, в которой за запросом q_1 следует запрос q_2 . Правило 3 отражает правило 1, указывая, что документ, выбранный в списке результатов для запроса q_2 , является более релевантным, чем предыдущие пропущенные документы в результатах для запроса q_1 , потому что оба запроса связаны одной целью. Аналогично правило 4 отражает правило 2.

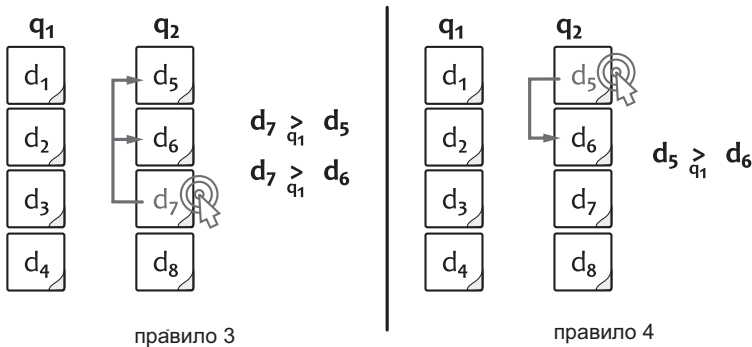


Рис. 4.31. Правила интерпретации неявной обратной связи для цепочки запросов

Последние два правила изображены на рис. 4.32. Они устанавливают отношения релевантности между документами из разных списков результатов в цепочке запросов. Правило 5 гласит, что виденные, но не выбранные документы в списке

результатов для запроса q_1 менее релевантны, чем выбранные документы, в наборе результатов для запроса q_2 . Это отношение релевантности устанавливается для предыдущего запроса. В соответствии с правилами 1 и 2 документы считаются виденными, если они находятся в списке выше выбранного или непосредственно под последним выбранным, как документ d_3 на рис. 4.32. Наконец, правило 6 гласит, что документы, выбранные в последнем списке результатов, релевантнее *двух* первых документов в первом списке. Это правило основано на предположении, что пользователь анализирует, по крайней мере, первые два результата в списке, прежде чем переформулировать запрос.

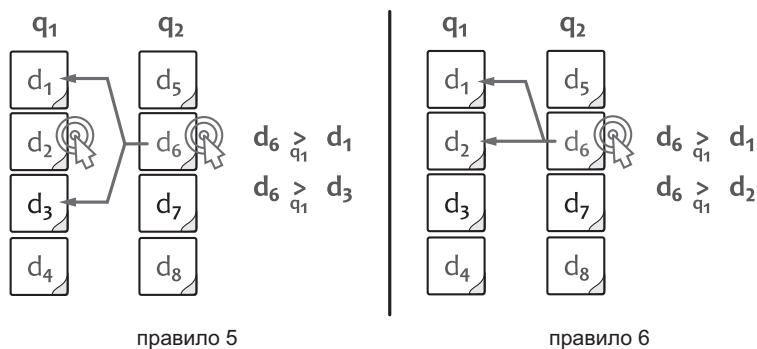


Рис. 4.32. Правила интерпретации неявной обратной связи для цепочки запросов (продолжение)

По каждой цепочке запросов в обучающем наборе данных одновременно оцениваются все шесть правил, чтобы получить отношения релевантности в форме

$$d_i >_q d_j, \quad (4.106)$$

означающие, что в отношении запроса q документ d_i релевантнее документа d_j . Эти правила можно непосредственно использовать в качестве обучающих меток в алгоритмах попарного обучения ранжированию, рассматривавшихся выше. Например, чтобы получить окончательную модель ранжирования, авторы только что описанной модели обратной связи использовали ее совместно с алгоритмом RankSVM [Radlinski and Joachims, 2005].

Неявная обратная связь несет важный сигнал, который можно использовать для автоматической настройки релевантности. Этот сигнал можно смешать с органическими оценками релевантности, полученными другими методами, такими как TF×IDF, чтобы недостатки в органическом ранжировании можно было исправить с помощью модели ранжирования, обученной на неявной обратной связи.

4.8. Архитектура служб поиска товаров

Завершим наше путешествие по методам поиска обзором обобщенной логической архитектуры службы поиска товаров, изображенной на рис. 4.33. Цель этого раздела — обобщить основные этапы обработки данных и запросов, затрагивавшиеся выше, не углубляясь в технические детали и детали реализации.

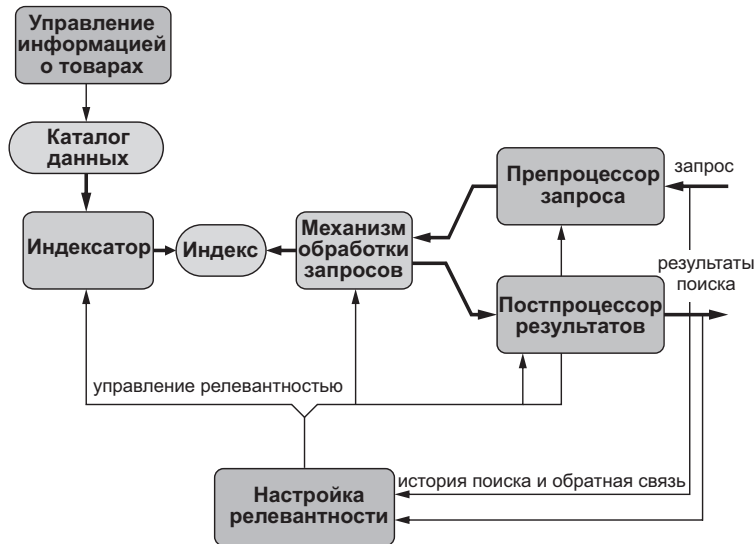


Рис. 4.33. Основные логические компоненты службы поиска товаров

Службу поиска можно рассматривать как базу данных, которая хранит каталог и реализует поиск в нем. Поскольку пользовательские запросы обычно обрабатываются в режиме реального времени, механизм поиска должен предварительно обработать и индексировать каталог с данными, чтобы обеспечить эффективную обработку запросов. В результате получается набор индексов, которые используются для сопоставления документов с запросом и оценки их соответствия. То есть простейшая служба поиска состоит из двух компонентов — индексатора и механизма обработки запросов. Обычно индексатор выполняет два основных этапа обработки данных: отображение и индексирование. Целью отображения является извлечение различных фрагментов входных данных и создание документов с четко определенными полями, значениями и иерархическими связями. Этап отображения обычно включает анализ содержимого, выполняющий лексемизацию, стемминг и другие преобразования для нормализации текста. Документы, созданные на этапе отображения, индексируются в структуры данных, обеспечивающие быструю обработку

запросов. Механизм обработки запросов реализует основные операции поиска, такие как сопоставление лексем, обработку логических запросов или вычисление оценки $TF \times IDF$ поверх индексов.

Многие методы поиска, такие как стемминг, n -граммы и расширение синонимов, требуют применения определенных преобразований как к полям документа, так и к запросам. На стороне документа эти преобразования часто выполняются индексатором и применяются к входным данным или документам до построения индекса, потому что индексы создаются на основе фактических лексем в документах и многие преобразования нельзя применить эффективно во время обработки запроса, если создавать индексы на основе необработанных и непреобразованных данных. На стороне запроса преобразования применяются во время его обработки, поэтому индексы и запросы приводятся к одной и той же нормальной форме. Нормализация — это не единственный тип преобразования, применяемый к запросам. Как говорилось выше, некоторые методы поиска, такие как контролируемое снижение точности, радикально преобразуют начальный запрос или генерируют несколько промежуточных запросов. Логика преобразования запроса часто реализуется в виде препроцессора запроса, который разбивает исходный запрос на основные примитивы, поддерживаемые ядром механизма обработки запросов. Индексатор и препроцессор должны реализовать одну и ту же логику преобразования и применять одинаковые алгоритмы стемминга, деления на n -граммы и семантического расширения.

Результатами механизма обработки запросов являются совпавшие документы и их оценки релевантности, рассчитанные с использованием базовых методов смешивания сигналов. Эти результаты могут преобразовываться постпроцессором результатов с применением правил группировки и продвижения, дополняющих или переопределяющих базовую оценку. Например, постпроцессор может реализовать правила повышения и понижения для подъема продвигаемых продуктов в результатах поиска.

Для эффективного поиска товаров требуется, по крайней мере, два процесса, действующих за кулисами и дополняющих только что описанный конвейер индексирования и обработки запросов. Первый — настройка релевантности, который управляет параметрами управления релевантностью во всех компонентах конвейера индексирования и обработки запросов и согласует их работу. Это может быть ручной процесс или компонент машинного обучения, анализирующий историю запросов и обратную связь от пользователей и оптимизирующий алгоритмы релевантности. Второй процесс — управление информацией о продукте (Product Information Management, PIM), задача которого заключается в очистке, подготовке и обогащении данных из каталога, загруженных в механизм поиска. Качество и полнота входных данных имеют решающее значение для качества по-

иска, потому что многообразие и точность генерируемых сигналов релевантности напрямую зависят от многообразия и точности атрибутов продукта. Например, базовых описаний продуктов может быть недостаточно для надлежащей обработки запросов, таких как *безглютеновая выпечка* или *платье с длинными рукавами*. Чтобы механизм поиска понимал такие запросы, товарам должны быть тщательно присвоены соответствующие характерные признаки, помогающие классифицировать их как продукты питания, одежду с рукавами и т. д. Общее количество полей, описывающих товары и индексируемых механизмом поиска, может достигать нескольких сотен. Получение этих метаданных и управление ими представляет собой серьезную проблему, поскольку разные части информации могут поступать от производителей, сторонних поставщиков данных или создаваться внутри организации. Маркетолог может упростить этот процесс, используя для создания метаданных и проверки их качества систему управления информацией о продукте и специализированные инструменты. Например, некоторые ретейлеры используют продвинутое инструменты распознавания изображений, чтобы получить или проверить определенные атрибуты продукта, такие как тип одежды или цвет, по изображениям продукта.

4.9. Итоги

- Цель служб поиска — получение предложений, соответствующих намерениям клиента, выраженным в запросе или в выбранных фильтрах. Службы поиска решают задачу открытия продукта, которую можно рассматривать как частный случай таргетирования.
- К основным компонентам среды поиска товаров относятся: пользовательский интерфейс для ввода запросов и отображения ранжированных документов, механизм поиска, обрабатывающий запросы и ранжирующий документы, и процесс настройки релевантности, оптимизирующий параметры управления релевантностью, которые определяют соответствие запросов и документов.
- Основными целями службы поиска являются релевантность, гибкость управления продвижением товаров и качество службы. Выгоды от поддержки службы поиска могут напрямую зависеть от релевантности и управления продвижением товаров.
- К основным метрикам релевантности относятся: точность/полнота и взвешенная накопленная релевантность в ранжированных результатах поиска.
- Средства управления релевантностью можно использовать для улучшения релевантности и достижения дополнительных бизнес-целей, таких как про-

движение определенных продуктов. К средствам управления релевантностью относятся такие методы, как повышение и понижение, фильтрация, фиксированные результаты, перенаправление на другие страницы и группировка продуктов.

- Метрики качества службы поиска включают: коэффициент конверсии, процент переходов, время, проведенное на странице с подробной информацией о продукте, частота изменения запроса, частота листания, коэффициент удержания и задержка поиска. Некоторые из этих метрик можно использовать в качестве целей при автоматической настройке релевантности.
- Поток обработки запросов можно рассматривать как многоступенчатый процесс, который разбивает документы и запросы на признаки, получает сигналы релевантности, сопоставляя эти два набора признаков, а затем использует сигналы для принятия решений о ранжировании.
- Набор основных методов поиска включает методы предварительной обработки текста (лексемизация, удаление стоп-слов, стемминг), сопоставления лексем и логического поиска. Самые простые методы оценки основаны на модели векторного пространства, которая представляет документы в виде векторов в линейном пространстве, где каждое измерение соответствует отдельному терму. Популярный метод оценки, $TF \times IDF$, уточняет базовую модель векторного пространства и производит оценку с использованием статистик частот термов.
- В реальной жизни службы поиска обычно используют документы с несколькими полями, которые можно оценить по отдельности и получить несколько сигналов релевантности. Эти сигналы можно смешивать, используя разнообразные методы проектирования сигналов.
- Методы сопоставления термов не способны определять семантические отношения, такие как синонимия и полисемия. Это ограничение устраняется методами семантического анализа. С точки зрения поиска большинство методов семантического анализа можно рассматривать как методы векторного представления слов, которые отображают слова, документы или запросы в векторы вещественных чисел с определенными семантическими свойствами. К ключевым методам семантического анализа можно отнести латентно-семантический анализ, вероятностное тематическое моделирование и контекстное векторное представление слов.
- Поиск товаров часто имеет дело со структурированными сущностями и конкретными требованиями к точности/полноте, которые невозможно удовлетворить с помощью обычных методов поиска. Промышленный опыт показывает, что хорошие результаты можно получить, используя логические методы, дающие высокую точность и низкую полноту.

- Настройка релевантности — это процесс оптимизации качества метрик поиска путем изменения параметров управления релевантностью. Эта задача тесно связана с классификацией и регрессией — для заданного запроса нужно предсказать ранг документа. Однако она отличается от стандартной классификации, то есть для ее решения применяются специализированные алгоритмы обучения ранжированию. Типичными признаками, которые используются в методах обучения ранжированию, являются статистики документов, статистики запросов, сигналы релевантности и неявная обратная связь от пользователя.
- К основным компонентам службы поиска товаров относятся: индексатор, основной механизм обработки запросов, препроцессор запросов, постпроцессор результатов и модули настройки релевантности.

5

Рекомендации

Разнообразие продуктов и услуг, предлагаемых клиентам, ограничено рядом факторов, включая затраты на производство и распространение. Продуктовый магазин может торговать лишь определенным количеством продуктов из-за ограниченности полочного пространства, радиостанция может втиснуть в свой ежедневный эфир не более определенного количества песен, а театр может поставить только ограниченное количество спектаклей. Продавец может увеличивать торговые площади и расширять ассортимент, но после определенного момента дополнительные доходы от расширения начинают уменьшаться из-за ограниченности общего спроса. С другой стороны, издержки, связанные с расширением, могут уменьшаться не так быстро, как доходы, или не уменьшаться вовсе, поэтому в какой-то момент расходы превысят доходы, что сделает дальнейшее расширение ассортимента экономически невыгодным. По этой причине продавцу, как правило, приходится ориентироваться на относительно популярные товары и предлагать лишь ограниченное разнообразие нишевых товаров.

Спрос на нишевые товары, однако, существует, что создает *длинный хвост*¹ в гистограмме популярности продукта, изображенной на рис. 5.1. На практике общий спрос на такие нишевые товары порой сопоставим с общим спросом на популярные продукты [Anderson, 2008]. Эти два общих требования соответствуют областям D_1 и D_2 под кривой спроса на рис. 5.1. Кроме того, товары с длинным хвостом часто могут быть высококачественными продуктами с более высокой доходностью, чем основные популярные товары, что делает еще более значительным вклад таких товаров в общую прибыль.

¹ Понятие «длинного хвоста» ассоциируется с обратной степенной зависимостью. В бизнесе концепция длинного хвоста была разработана Крисом Андерсоном для описания явления больших суммарных продаж товаров, ставших в свое время классикой, по сравнению с товарами, которые в настоящее время считаются модными. — *Примеч. ред.*

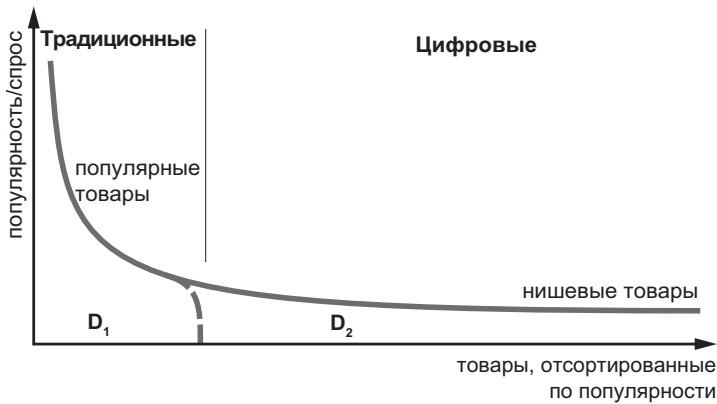


Рис. 5.1. Типичная гистограмма популярности продуктов

Как уже упоминалось, многие традиционные каналы распространения, такие как обычные магазины, театры и радиостанции, имеют ограниченную возможность удовлетворить спрос с длинным хвостом из-за затрат на распространение. Однако развитие цифровых каналов резко изменило ситуацию. Во-первых, новые медиа-каналы практически исключили расходы на распространение цифрового контента и дали возможность создавать онлайн-службы с практически бесконечным ассортиментом. Например, онлайн-служба по распространению видеофильмов может предложить практически неограниченное и постоянно расширяющееся разнообразие видеофильмов, включая голливудские фильмы, телесериалы со всего мира и любительские фильмы. Во-вторых, цифровые каналы дали ретейлерам и производителям нецифровых товаров возможность изменить структуру издержек на распространение и охватить большее число потребителей. Клиентская база обычного магазина ограничена людьми, живущими неподалеку или часто посещающими этот район, поэтому разнообразие спроса также ограничено. Интернет-магазин, охватывающий всю страну или весь мир, имеет дело с гораздо более разнообразным спросом. Благодаря этому появились и преуспели продавцы с гигантским ассортиментом, такие как Amazon. Впечатляющее увеличение ассортимента с сильным акцентом на нишевые продукты усложняет открытие продуктов с использованием старых подходов, потому что средний клиент способен просмотреть только небольшую часть доступных предложений, которых могут быть миллионы. Эта потребность в мощных услугах по открытию продуктов стала одной из основных движущих сил развития рекомендательных систем.

Сервисы рекомендаций, в отличие от поисковых служб, нацелены на предоставление клиенту релевантных предложений, когда цель поиска не выражена или выражена нечетко. Иногда цель поиска трудно выразить явно из-за сложностей,

связанных с формализацией требуемых свойств продукта. Например, цель поиска для клиента, ищущего музыку, может определяться его личными вкусами, но такие намерения порой трудно перевести в формальные критерии. В других случаях клиент может не знать о некоторых видах или категориях продуктов, или просто не осознавать, или сомневаться в собственных потребностях. Фотограф-любитель, например, может не понимать, что для достижения наилучших результатов в режиме фотосъемки, который его интересует, требуются специальные объективы. Следовательно, в отличие от служб поиска, которые получают запрос и могут оценивать соответствие продуктов этому запросу, рекомендательная система должна угадывать намерения покупателя, опираясь на косвенную информацию, такую как рейтинги продуктов и история покупок клиента. На основе этой информации можно вычислять различные метрики сходства и использовать их как альтернативу оценкам сходства запрос/продукт, вычисляемым службами поиска. В частности, рекомендательная система может использовать следующие сходства.

СХОДСТВО С ДРУГИМИ ПОЛЬЗОВАТЕЛЯМИ. Цель покупки данного клиента можно вывести из поведения похожих клиентов в прошлом. Этот подход напоминает моделирование методом аналогии, которое мы рассмотрели выше.

СХОДСТВО С ДРУГИМИ ПРОДУКТАМИ. Для определения групп продуктов и категорий, наиболее релевантных для данного клиента, можно использовать поиск и покупки товаров в прошлом и рекомендовать похожие продукты.

СХОДСТВО С ДРУГИМИ КОНТЕКСТАМИ. Точность рекомендаций можно повысить, используя не только характеристики клиента и продукта, но также контекстную информацию, которая несет дополнительные сигналы о намерении совершить покупку. Например, продавец модной одежды может рекомендовать очень разные продукты для одного и того же клиента в зависимости от сезона.

Алгоритмически методы выработки рекомендаций имеют много общего с методами поиска и так же используют предиктивные методы, которые мы использовали ранее для таргетирования рекламных акций. Остальная часть этой главы представляет систематическое описание рекомендательных систем, начиная с окружающей среды и экономических целей, а затем углубляется в различные методы выработки рекомендаций.

5.1. Среда

По своим основным параметрам службы рекомендаций похожи на службы поиска. Основная цель рекомендательной системы, как и службы поиска, — представить клиенту ранжированный список рекомендуемых товаров. Эти рекомендации могут

передаваться через различные маркетинговые каналы. Мы будем считать, что рекомендации запрашиваются по каналу, действующему в режиме реального времени, что характерно для веб-сайтов и мобильных приложений, хотя некоторые каналы, такие как электронная почта, могут иметь более мягкие требования и позволяют рассчитывать рекомендации в автономном режиме. Однако перечень основных элементов рекомендательной системы, изображенных на рис. 5.2, отличается от перечня элементов служб поиска и включает:

- Подавляющее большинство методов выработки рекомендаций предполагает, что для элементов каталога доступны *рейтинги по оценкам клиентов*. Рейтинги могут явно определяться клиентами или вычисляться на основе поведенческих данных, таких как покупки и истории просмотров. Каждое значение рейтинга представляет обратную связь от определенного клиента по определенному товару и измеряется по определенной шкале. Обратите внимание, что клиент может оценить любые товары в каталоге, а не только рекомендованные. Иначе говоря, рейтинги отражают мнение клиентов о товарах в каталоге, а не рекомендации. Более подробно о рейтингах и их свойствах мы поговорим в следующем разделе.

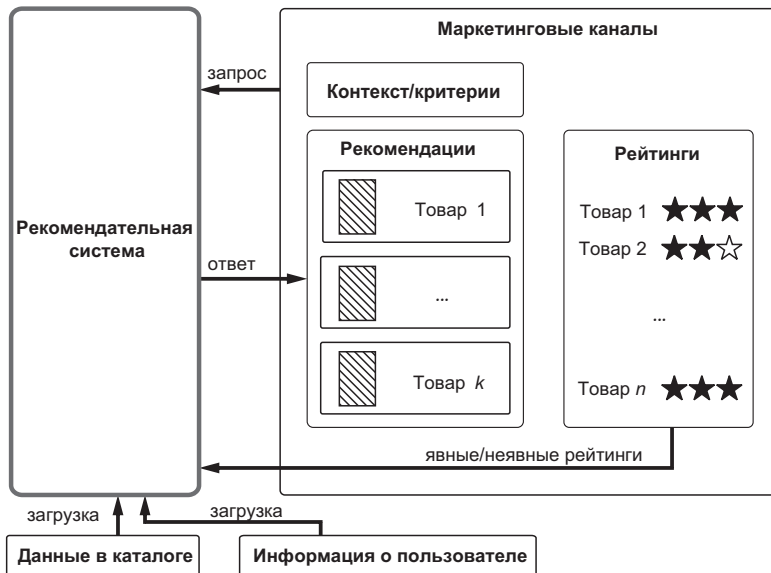


Рис. 5.2. Основные элементы среды службы рекомендаций

- Некоторые методы выработки рекомендаций опираются на содержимое и данные в каталоге для вычисления сходства товаров на основе их характеристик.

Для этого, по аналогии с поисковыми системами, рекомендательные системы должны интегрироваться с источниками данных, такими как система управления информацией о продуктах.

- Некоторые рекомендательные системы могут использовать дополнительные данные о пользователях, такие как истории онлайн-заказов или покупки в магазинах. Эту информацию следует отличать от информации о рейтингах: рейтинги лишь описывают взаимодействия пользователь/товар, тогда как внешние данные, такие как история заказов, могут дать дополнительные сведения о группировке взаимодействий в заказы и т. д.
- Запросы на получение рекомендаций и рейтинги клиентов могут дополняться контекстной информацией, например временем, местоположением или типом маркетингового канала. Рекомендательная система может использовать контекстные данные для повышения релевантности рекомендаций. Например, время установления рейтинга можно использовать для учета сезонных и временных тенденций во вкусах клиентов. Запрос на рекомендации или контекст также может включать явные критерии или предпочтения клиента, которые также желательно использовать для уточнения рекомендаций. Например, рекомендации средств по уходу за кожей можно корректировать в зависимости от типа кожи (нормальная, сухая, жирная и т. д.), указанного в предпочтениях клиента.

Разные семейства методов выработки рекомендаций могут использовать разные подмножества данных, и их достоинства и недостатки в значительной степени определяются диапазоном доступных данных.

5.1.1. Свойства рейтингов клиентов

Рейтинги клиентов часто считаются наиболее важным источником информации для выработки рекомендаций, поэтому мы должны внимательно изучить, как определяются рейтинги, и познакомиться с типичными свойствами рейтингов.

Обычно рейтинги представлены в виде матрицы, где строки соответствуют пользователям, а столбцы — товарам. Обозначим матрицу рейтингов как $\mathbf{R} = (r_{ij})$, где r_{ij} — это рейтинг, присвоенный пользователем i товару j . В рекомендательной системе, отслеживающей m пользователей и содержащей каталог из n элементов, матрица \mathbf{R} имеет размер $m \times n$. На практике матрица рейтингов почти всегда неполная — рейтинги известны только для подмножества пар пользователь/товар, а остальные элементы отсутствуют (не указаны). Чаще рейтинги представлены числовыми значениями, которые могут определяться по-разному, в зависимости от сферы бизнеса, маркетингового канала и источника данных. Выделим следующие два варианта:

ПОРЯДКОВЫЙ. Интерфейс рекомендательной системы часто позволяет пользователю выразить свои предпочтения, выбирая рейтинги из дискретного набора чисел (например, 1, 2 или 3 звезды) или непрерывного диапазона (например, от -5 до $+5$). Нередки рейтинги в виде оценок с двумя (например, «нравится» или «не нравится») или большим (например, «плохо», «хорошо» или «отлично») числом категорий, которые затем отображаются в дискретные числовые значения. Неявную обратную связь тоже можно выразить в виде порядковых значений, но такие значения обычно отражают *уверенность*, а не *предпочтения*. Например, неявный рейтинг может указывать, как часто пользователь покупает определенный товар или сколько времени клиент провел на веб-странице с описанием товара.

УНАРНЫЙ. Во многих случаях матрица рейтингов отражает не уровень близости между пользователем и товаром, а лишь сам факт взаимодействия. Например, многие интерфейсы имеют только одну кнопку **Нравится**, поэтому пользователю остается лишь либо подтвердить, что ему нравится товар, либо вообще не оставлять никакой обратной связи. Другим типичным примером унарных рейтингов является неявная обратная связь, регистрирующая взаимодействие между пользователями и товарами, но не фиксирующая таких деталей, как количество покупок, хотя можно смело утверждать, что простые количественные свойства этой обратной связи имеют большое значение [Hu et al., 2008]. Элементы унарной рейтинговой матрицы могут принимать только два значения — «единица» или «не указано».

Значения рейтинга в матрице **R** часто дополняются контекстной информацией, такой как дата и время оценки рейтинга, маркетинговый канал, используемый клиентом установки рейтинга, и т. д. Эту информацию можно использовать в рекомендательной системе, чтобы определить наиболее релевантные значения рейтинга для данного контекста.

Важно отметить, что значения рейтингов в матрице с порядковыми оценками также несут *неявную обратную связь*. Дело в том, что пользователи более склонны оценивать товары, которые им нравятся, и избегать товаров, которые их не привлекают. Например, пользователь может полностью избегать музыки определенных жанров или товаров из определенных категорий. Поэтому кроме самих оценок важно также учитывать, какие товары оцениваются. Другими словами, распределение рейтингов для случайных товаров, вероятно, будет отличаться от распределения рейтингов товаров, выбранных пользователем. Строго говоря, это означает, что рекомендательная система не должна полагаться на предположение, что распределение наблюдаемых рейтингов является репрезентативным для распределения отсутствующих рейтингов [Devooght et al., 2015]. Как будет показано ниже, некоторые передовые методы выработки рекомендаций принимают это во внимание и выводят неявную обратную связь из матрицы рейтингов. В более общем случае

рекомендательный алгоритм может вовлекать в расчеты две отдельные матрицы рейтингов — для явной и неявной обратной связи.

Второе важное свойство матрицы рейтингов — *разреженность*. Матрица рейтингов по своей сути является разреженной, так как любой отдельный пользователь взаимодействует только с небольшой долей доступных товаров, соответственно каждая строка матрицы содержит только несколько известных рейтингов, а все остальные значения отсутствуют. Более того, распределение известных рейтингов обычно демонстрирует свойство длинного хвоста, о котором говорилось выше. Это означает, что непропорционально большое количество известных рейтингов соответствует лишь нескольким наиболее популярным товарам, тогда как рейтинги нишевых продуктов особенно немногочисленны. Это свойство можно проиллюстрировать на примере известного набора данных, опубликованных службой потокового видео Netflix, с рейтингами фильмов, выставленными примерно 500 000 подписчиками. Как оказывается, около 33 % рейтингов охватывают только 1,7 % самых популярных позиций [Cremonesi et al., 2010]. Свойство длинного хвоста представляет большую сложность при разработке и оценке рекомендательных систем, потому что алгоритмы выбора и метрики качества рекомендаций имеют тенденцию смещаться в сторону популярных элементов, что снижает качество рекомендаций для нишевых продуктов.

5.2. Бизнес-цели

Ключевые бизнес-цели сервисов рекомендаций тесно связаны с целями поиска товаров. Основные соображения, рассмотренные выше в контексте поиска товаров, такие как релевантность и управление продвижением, в целом применимы к рекомендательным системам. Главное отличие заключается в том, что цель поиска не выражена явно и вообще может отсутствовать. Это требует расширения основной цели представления релевантных результатов, потому что понятие релевантности становится более шатким по мере того, как цель поиска теряет свою направленность. Следовательно, стандартный набор целей рекомендательной системы часто определяется следующим образом.

РЕЛЕВАНТНОСТЬ. Рекомендации, предложенные пользователю, должны быть релевантными в том смысле, что пользователь должен иметь высокую склонность приобретать рекомендуемые товары и высоко оценивать их.

НОВИЗНА. Рекомендательная система не обрабатывает явный поисковый запрос, а скорее консультирует пользователей по доступным вариантам. Следовательно, рекомендательная система должна предоставлять варианты, пока не известные пользователям; иначе рекомендации могут восприниматься как

тривиальные и бесполезные. Типичным примером этой проблемы является рекомендация похожих популярных товаров, о которых пользователь наверняка будет знать. Например, пользователь, купивший одну из книг о Гарри Поттере, может получить рекомендацию купить другие книги из этой серии, а не другие книги того же жанра.

СЕРЕНДИПНОСТЬ. Рекомендации могут помочь пользователю обнаружить неожиданные и удивительные продукты, которые также являются новыми. Например, пользователю, покупающему книги о машинном обучении, можно предложить купить другие книги по этой теме — и хотя некоторые из них могут оказаться новыми для пользователя, их едва ли можно рассматривать как приятную неожиданность. С другой стороны, рекомендательная система может попытаться угадать бизнес-сферу, интересующую пользователя, и рекомендовать книгу по аналитическим методам для конкретной области, например, по анализу поведения клиентов или торговым моделям, которые могут оказаться приятной неожиданностью. Еще более интересными примерами неожиданных рекомендаций могут служить рекомендации с товарами из совершенно других категорий. Например, пользователю, покупающему книги о Гарри Поттере, можно рекомендовать посетить парк развлечений с соответствующей достопримечательностью или пользователям, изучающим европейские средневековые эпосы, такие как «Беовульф» и «Песнь о Роланде», можно рекомендовать билеты в оперу. Приятные и неожиданные рекомендации не только повышают качество службы и коэффициент конверсии, но также могут помочь создать прочную основу для долгосрочных отношений с клиентом.

РАЗНООБРАЗИЕ. Наконец, список предлагаемых рекомендаций должен быть разнообразным, чтобы увеличить вероятность конверсии пользователя. Список рекомендаций, состоящий из очень похожих элементов, даже если они релевантны, новы и неожиданны, может оказаться неоптимальным.

По аналогии со службами поиска, общую прибыль от рекомендательной системы можно определить с точки зрения доходности продукта и проданного количества:

$$\text{Profit} = \sum_{\text{products}} \text{Quantity sold}_{\text{product}} \times \text{Margin}_{\text{product}}. \quad (5.1)$$

Релевантность, новизна, серендипность и разнообразие — все эти цели направлены на повышение коэффициента конверсии и, как следствие, объема продаж. Рекомендации, выработанные в соответствии с этими целями, можно повторно ранжировать для достижения дополнительных маркетинговых целей, таких как продвижение высокодоходных или сезонных продуктов, чтобы увеличить член уравнения, определяющий доходность от продажи продукта.

5.3. Оценка качества

Наш следующий шаг — разработка количественных показателей для оценки качества рекомендательных систем с учетом целей, перечисленных в предыдущем разделе.

Качество результатов поиска, как правило, можно оценить с использованием экспертных суждений о релевантности элементов в контексте данного запроса, но этот подход имеет ограниченную применимость для рекомендаций, поскольку контекст обычно включает данные профиля пользователя и, следовательно, является уникальным для каждого пользователя. Это затрудняет или делает невозможным оценку качества рекомендаций вручную для каждого возможного контекста. С другой стороны, матрица рейтингов уже содержит экспертные оценки, представленные пользователями в их собственных персонализированных контекстах. Следовательно, проблему рекомендаций можно рассматривать как проблему прогнозирования рейтингов (как рекомендовать элементы с самыми высокими прогнозируемыми рейтингами для данного пользователя), а качество рекомендаций можно измерить сравнением прогнозируемых и фактических рейтингов, известных из матрицы рейтингов. С этой точки зрения задача рекомендаций очень близка к классификации или регрессии.

Напомню, что задачу классификации/регрессии можно определить с помощью матрицы, в которой каждая строка представляет точку данных, а столбцы — признаки или отклики. Предиктивная модель обучается на данных как с признаками, так и с откликами, а затем используется для прогнозирования ответов на основе признаков. Это проиллюстрировано в примере 5.2, где точки данных 1–3 используются для обучения, а фактическое предсказание выполняется для точек 4–6. Кроме того, для обучения и настройки модели точки данных с известными откликами обычно делятся на обучающий, проверочный и контрольный наборы данных. Первоначально модель конструируется с использованием обучающего набора. Затем прогнозируются отклики для проверочного набора, сравниваются с фактическими значениями и оценивается качество модели. По результатам оценки модель может быть перестроена с другими параметрами, вновь обучена на обучающих данных и оценена повторно. Набор контрольных данных используется для окончательной оценки качества модели в самом конце процесса.

$$\begin{array}{l}
 \text{Признак 1} \quad \text{Признак 2} \quad \text{Признак 3} \quad \text{Отклик} \\
 \begin{array}{l}
 \text{Точка данных 1} \\
 \text{Точка данных 2} \\
 \text{Точка данных 3} \\
 \text{Точка данных 4} \\
 \text{Точка данных 5} \\
 \text{Точка данных 6}
 \end{array}
 \begin{bmatrix}
 x_{11} & x_{12} & x_{13} & y_1 \\
 x_{21} & x_{22} & x_{23} & y_2 \\
 x_{31} & x_{32} & x_{33} & y_3 \\
 x_{41} & x_{42} & x_{43} & - \\
 x_{51} & x_{52} & x_{53} & - \\
 x_{61} & x_{62} & x_{63} & -
 \end{bmatrix}
 \end{array} \quad (5.2)$$

Ключевое отличие задачи прогнозирования рейтинга состоит в том, что в матрице рейтингов отсутствуют признаки и отклики. Известные и неизвестные рейтинги смешиваются без какой-либо конкретной структуры, как показано в примере 5.3, и цель состоит в том, чтобы обучить модель предсказывать неизвестные рейтинги по известным. Эта задача заполнения недостающих элементов частично наблюдаемой матрицы известна как задача *восстановления матрицы*.

$$\begin{array}{l}
 \text{Пользователь 1} \\
 \text{Пользователь 2} \\
 \text{Пользователь 3} \\
 \text{Пользователь 4}
 \end{array}
 \begin{array}{c}
 \text{Элемент 1} \quad \text{Элемент 2} \quad \text{Элемент 3} \quad \text{Элемент 4} \\
 \left[\begin{array}{cccc}
 r_{11} & - & r_{13} & - \\
 - & r_{22} & - & r_{24} \\
 - & r_{32} & r_{33} & - \\
 r_{41} & - & r_{43} & r_{44}
 \end{array} \right]
 \end{array}
 \quad (5.3)$$

Так же как в случае со стандартной задачей классификации, мы должны разделить имеющиеся данные на обучающий, проверочный и контрольный наборы, чтобы обучить модель и оценить ее качество. В задачах классификации это делается по строкам. Например, доступные данные в матрице 5.2 можно разделить на три набора, назначив первую строку обучающим набором, вторую — проверочным набором и третью — контрольным. Этот подход плохо подходит для решения задачи восстановления матрицы, поскольку подразумевает, что модель обучается на одном наборе пользователей и оценивается на другом, что на самом деле не так. Поэтому матрица рейтингов обычно разбивается на наборы по элементам. То есть из исходной матрицы удаляется определенная часть известных рейтингов, из оставшихся элементов формируется обучающая матрица, а удаленные рейтинги помещаются в проверочный и контрольный наборы, которые затем используются для оценки качества прогнозирования.

Интерпретируя задачу получения рекомендаций как прогнозирование рейтинга, можно определить несколько метрик качества, связанных с бизнес-целями. В следующих разделах мы разработаем эти метрики.

5.3.1. Точность прогнозирования

Точность прогнозирования рейтингов можно рассматривать как меру релевантности, потому что она количественно определяет, насколько хорошо рекомендательная система предсказывает полезность для пользователя, оцениваемую самими пользователями. Для оценки точности у нас имеется широкий выбор метрик, используемых в машинном обучении и информационном поиске, включая метрики качества поиска, рассмотренные выше.

Первое семейство метрик, которые мы рассмотрим, — метрики точности прогнозирования, широко используемые для оценки методов классификации и регрессии.

Обозначим множество наблюдаемых рейтингов $r_{ij} \in \mathbf{R}$ как R , а проверочное подмножество, используемое для оценки точности, как $T \subset R$. Для каждого рейтинга в T рекомендательный алгоритм вычисляет оценку \hat{r}_{ij} , то есть ошибку прогнозирования можно определить так:

$$e_{ij} = \hat{r}_{ij} - r_{ij}. \quad (5.4)$$

Общее качество прогнозирования рейтингов можно получить усреднением ошибок поточечных прогнозов. Есть несколько способов определения этой усредненной метрики. Первый — среднеквадратическая ошибка (Mean Squared Error, MSE), определяемая как

$$\text{MSE} = \frac{1}{|T|} \sum_{(u,j) \in T} e_{ij}^2. \quad (5.5)$$

Метрика MSE не всегда удобна, потому что она оперирует квадратами ошибок, которые нельзя сравнивать с исходными рейтингами непосредственно. Эту проблему можно устранить, взяв корень от среднеквадратической ошибки (Root Mean Squared Error, RMSE), — эта величина измеряется в тех же единицах, что и исходные рейтинги:

$$\text{RMSE} = \sqrt{\text{MSE}}. \quad (5.6)$$

Кроме того, величину RMSE можно нормализовать, приведя ее к диапазону (0, 1), и получить нормализованную метрику RMSE (Normalized RMSE, NRMSE):

$$\text{NRMSE} = \frac{\text{RMSE}}{r_{\max} - r_{\min}}. \quad (5.7)$$

Благодаря своей простоте метрика RMSE и ее вариации широко используются на практике для оценки рекомендательных систем. Однако RMSE и другие подобные метрики оценки поточечной точности имеют несколько существенных недостатков:

- Как отмечалось выше, рейтинги обычно подчиняются распределению с длинным хвостом, более плотным для популярных элементов и разреженным для элементов из длинного хвоста. Это усложняет прогнозирование рейтингов для элементов из длинного хвоста, по сравнению с популярными элементами, и, как следствие, к различной точности прогнозирования для этих двух групп элементов. RMSE не различает эти две группы и просто берет среднее значение, поэтому плохая точность для элементов из длинного хвоста может уравниваться высокой точностью для популярных предметов. Чтобы измерить и устранить эту дилемму, можно рассчитать RMSE отдельно для раз-

ных групп элементов или добавить весовые коэффициенты в уравнение 5.5 и таким образом учесть особенности элементов или другие аспекты.

- Цель рекомендательной системы — предсказать по историческим данным, как пользователь оценит элемент в будущем. Вкусы и интересы пользователей могут меняться с течением времени, поэтому система должна распознавать такие временные тенденции, чтобы точно прогнозировать будущее поведение. RMSE непосредственно не учитывает этот аспект. Проблему, однако, можно решить, спроектировав надлежащий контрольный набор T . Чтобы проверить возможность прогнозирования рейтингов в будущем, можно выбрать контрольный набор T из рейтингов R не случайно, а так, чтобы обучающий набор содержал более старые рейтинги, а контрольный набор T — более свежие. Этот подход несколько противоречит стандартной методологии оценки моделей, поскольку обучающие и контрольные наборы, сконструированные таким способом, имеют разные распределения, но это действенный практический метод, который использовался, например, в Netflix Prize, открытом конкурсе на лучший алгоритм совместной фильтрации, проводимом Netflix, онлайн-службой потокового видео, в 2006–2009 годах [Aggarwal, 2016].

Это соображение особенно очевидно в случае рейтингов, полученных на основе неявной обратной связи. Например, если рейтинги определяются по событиям покупки, прогнозирование будущих рейтингов фактически означает прогнозирование будущих покупок.

- Система рекомендаций передает пользователю ранжированный список рекомендаций, обычно ограниченный K верхними элементами. RMSE не учитывает ранжирование и одинаково штрафует за ошибки прогнозирования элементы как в верхней, так и в нижней части списка. Можно утверждать, что алгоритмы с очень малой разницей в RMSE могут иметь большую разницу в первых K элементах в своих списках [Koren, 2007].

5.3.2. Точность ранжирования

Чтобы измерить качество K лучших рекомендаций, можно использовать богатый набор методов и метрик, разработанных для служб поиска. Прежде всего следует отметить, что понятия точности и полноты прямо связаны с задачей выбора K верхних рекомендаций. Если положить, что I_u — подмножество элементов в контрольном наборе T , положительно оцененных (например, купленных) пользователем u , и $Y_u(K)$ — список лучших K элементов, рекомендованных этому пользователю, мы можем определить метрики точности и полноты как функции от K :

$$\text{precision}(K) = \frac{|Y_u(K) \cap I_u|}{|Y_u(K)|}, \quad (5.8)$$

$$\text{recall}(K) = \frac{|Y_u(K) \cap I_u|}{|I_u|}. \quad (5.9)$$

С помощью этих двух метрик можно измерить качество рекомендательного алгоритма для любого заданного K . Точность (*precision*) — это процент релевантных рекомендаций в списке, а полнота (*recall*) — процент элементов, выбранных из набора доступных релевантных элементов. Два рекомендательных алгоритма можно сравнить в терминах точности и полноты, усредненных по пользователям, по аналогии с методами поиска. Число рекомендаций в списке, однако, является важным параметром, определяющим компромисс между точностью и полнотой. В коротких рекомендательных списках, как правило, отсутствуют релевантные элементы, тогда как в длинных списках обычно отмечается высокий процент нерелевантных. Этот компромисс можно визуальнo представить в виде кривой точность/полнота, о которой мы также говорили выше в контексте служб поиска. Кривая отображает значения точности и полноты для разных значений K и позволяет увидеть диапазон компромиссов между точностью и полнотой, достижимых алгоритмом рекомендаций.

Недостаток кривой точность/полнота состоит в том, что она не дает единой числовой метрики, обобщающей качество метода рекомендаций. К счастью, мы уже ввели ряд метрик качества ранжирования, которые можно адаптировать для получения именно такого обобщенного представления. Например, для рекомендаций можно адаптировать взвешенную накопленную релевантность (Discounted Cumulative Gain, DCG), используя известные рейтинги в качестве степеней релевантности. Напомню, что мы определили DCG для списка K элементов следующим образом:

$$\text{DCG} = \sum_{i=1}^K \frac{2^{g_i} - 1}{\log_2(i+1)}, \quad (5.10)$$

где g_i — степень релевантности i -го элемента списка. Если контрольный набор T содержит оценки m пользователей, мы можем определить общую DCG как среднее значение оценок DCG списков рекомендаций для отдельных пользователей:

$$\text{DCG} = \frac{1}{m} \sum_{u=1}^m \sum_{\substack{i \in I_u \\ R_{ui} \leq K}} \frac{2^{r_{ui}} - 1}{\log_2(R_{ui} + 1)}, \quad (5.11)$$

где I_u — подмножество элементов в контрольном наборе T , положительно оцененных пользователем u , R_{ui} — ранг элемента в списке рекомендаций для пользователя u , и r_{ui} — рейтинг из множества T , присвоенный пользователем u элементу i , который используется как аппроксимация степени релевантности g_i из уравнения 5.10. Обратите внимание, что внутренняя сумма в уравнении 5.11 просто перебирает

верхние K рекомендаций с тестовыми рейтингами, известными для данного пользователя. Другие стандартные метрики из информационного поиска, такие как нормализованная DCG (NDCG) и усредненная средняя точность (MAP), можно переформулировать аналогичным образом.

5.3.3. Новизна

Рекомендации считаются новыми, если пользователь не был знаком с рекомендуемыми элементами на момент рекомендации. Эта информация отсутствует в матрице рейтингов, поэтому ее следует собирать с помощью тестов и опросов пользователей или как-то иначе выводить из матрицы рейтингов. Поскольку тестирование и опросы, как правило, требуют много времени и ресурсов, можно попытаться разработать метрику новизны на основе матрицы рейтингов, сделав определенные предположения. Один из возможных подходов — обучение алгоритма рекомендаций на более старых рейтингах и его оценка с использованием более поздних рейтингов, как показано на рис. 5.3.

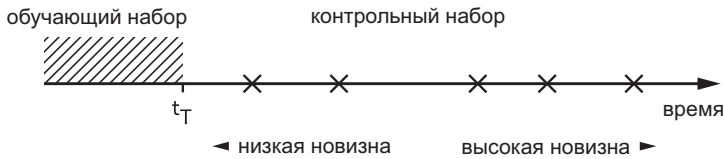


Рис. 5.3. Оценка новизны рекомендаций

Предполагается, что алгоритм рекомендаций, прогнозирующий элементы, ранжированные или купленные сразу за границей t_T обучающего набора, обеспечивает более низкую новизну, чем алгоритм, прогнозирующий элементы, ранжированные или купленные в более отдаленном будущем, потому что недавно приобретенные товары, вероятнее всего, уже известны пользователю. Следовательно, метрика новизны может использовать взвешенные по времени оценки для подъема точных долговременных прогнозов и понижения краткосрочных.

5.3.4. Серендипность

Приятная неожиданность, или серендипность — это мера привлекательности и удивительности рекомендаций для пользователя [Herlocker et al., 2004]. Оценка, насколько рекомендации являются приятной неожиданностью, является еще более сложной задачей, чем оценка новизны, потому что это свойство очень субъективно и информация обратной связи обычно не дает никакого намека на уровень прият-

ности и неожиданности. Однако для его измерения можно разработать эвристические методологии. Один из возможных подходов — сравнение рекомендаций, полученных с помощью оцениваемого алгоритма, с рекомендациями, полученными с помощью некоторого базового алгоритма, который, как известно, предлагает тривиальные рекомендации, не являющиеся приятной неожиданностью [Ge et al., 2010]. Если множество элементов, рекомендованных пользователю оцениваемым алгоритмом, обозначить как Y_u , а множество элементов, рекомендованных базовым алгоритмом, обозначить как Y_u^0 , тогда меру приятной неожиданности можно выразить как

$$\text{serendipity} = \frac{1}{m \cdot K} \sum_{u=1}^m \sum_{i \in I_u} \mathbb{I}(i \in (Y_u \setminus Y_u^0)), \quad (5.12)$$

где m — число пользователей, I_u — множество элементов в контрольном наборе, положительно оцененных пользователем, K — количество рекомендаций в списке, и $\mathbb{I}(\cdot)$ — индикаторная функция, возвращающая истинное значение, если элемент принадлежит множеству Y_u и не принадлежит множеству Y_u^0 . Другими словами, эта мера приятной неожиданности оценивает рекомендательную систему на основе доли нетривиальных и релевантных элементов в списке рекомендаций.

5.3.5. Разнообразие

Разнообразие — это способность рекомендательной системы составить список рекомендаций, состоящий из разнородных элементов. Большое разнообразие обычно предпочтительнее, потому что увеличивает вероятность, что список будет содержать хотя бы несколько элементов, релевантных для пользователя. Высокое разнообразие также может быть предпочтительным с точки зрения продвижения товаров, поскольку способствует перекрестным продажам и широкому охвату каталога.

Для измерения разнообразия можно использовать метрики сходства, разработанные ранее для служб поиска. Например, можно вычислить косинусные расстояния между описаниями для всех пар элементов в списке рекомендаций и оценить разнообразие как обратную величину среднего расстояния.

5.3.6. Охват

Цель рекомендательной системы — предсказать недостающие рейтинги в матрице. Как будет показано ниже, многие алгоритмы рекомендаций полагаются на сходства элемент–элемент или пользователь–пользователь, вычисленные на основе матрицы рейтингов, поэтому предсказать рейтинги для элементов или пользователей, не

имеющих большого числа общих рейтингов с другими элементами и пользователями, может оказаться непростой задачей. В связи с этим важно измерить охват, обеспечиваемый рекомендательной системой, то есть процент пользователей или элементов, по которым система может давать рекомендации. В некоторых случаях этот процент можно оценить по требованиям, предъявляемым алгоритмом рекомендаций. Например, алгоритм может потребовать наличия у пользователя не менее пяти оценок, чтобы иметь право на получение рекомендаций. В общем случае рекомендательная система может предсказать рейтинг для любой пары пользователей и элементов, просто используя значение по умолчанию или случайное значение. Это означает, что особый интерес может представлять возможность оценки надежности прогнозируемых рейтингов (вероятность, что оценочное значение является точным) и компромисса между охватом и точностью путем исключения из оценки точности пользователей или элементов с наименее надежными рейтингами.

Альтернативой охвату является так называемый *охват каталога* [Ge et al., 2010]. Охват каталога определяется как процент элементов, *фактически* рекомендованных пользователям. Проблема заключается в том, что рекомендательная система может оценить рейтинги широкого спектра товаров, но верхние K рекомендаций в списках, представленных пользователям, могут по-прежнему включать почти одинаковые рекомендации, что фактически приравнивается к плохому охвату с точки зрения продвижения товаров. Метрику охвата каталога можно определить как процент элементов, отображаемых хотя бы в одном списке рекомендаций:

$$\text{catalog coverage} = \frac{1}{n} \left| \bigcup_{u=1}^m Y_u \right|, \quad (5.13)$$

где n — общее количество элементов в каталоге. Для оценки числа охватываемых позиций уравнение 5.13 использует объединение списков рекомендаций для всех пользователей. Альтернативный подход заключается в подсчете числа различных элементов, рекомендуемых в большом числе реальных сеансов пользователей.

5.3.7. Роль экспериментов

Метрики, описанные выше, помогают измерить качество рекомендаций с нескольких важных точек зрения. Однако конечной целью рекомендательной системы является повышение доходов и процента конверсии. Представленные метрики обеспечивают прочную основу для оценки качества, тем не менее они не имеют тесной связи с показателями финансовой деятельности. Эту связь можно установить путем практических экспериментов, многомерного тестирования и измерения подъема.

5.4. Обзор методов рекомендаций

До сих пор мы описывали среду и источники данных, с которыми интегрируется рекомендательная система, ее бизнес-цели и метрики, которые можно использовать для оценки качества рекомендаций. Это обеспечивает достаточно прочную основу для разработки алгоритмов рекомендаций. К этой задаче можно подойти с разных точек зрения, и существует несколько семейств методов, отличающихся источниками данных, используемыми для составления рекомендаций (матрица рейтингов, каталог или контекстная информация), и типами моделей прогнозирования рейтингов. В оставшейся части главы мы последовательно обсудим все основные категории алгоритмов рекомендаций, но прежде чем углубляться в детали, будет полезно кратко познакомиться с классификацией методов и сделать несколько общих замечаний.

Методы рекомендаций можно классифицировать по-разному, в зависимости от точки зрения. С позиции алгоритмического и информационного поиска методы рекомендаций классифицируются в первую очередь по типу предиктивной модели и входным данным. Соответствующая иерархия показана на рис. 5.4. Исторически сложилось так, что двумя основными семействами методов рекомендаций являются фильтрация по содержимому и совместная фильтрация. Фильтрация по содержимому в основном опирается на содержимое — текстовые описания элементов, а совместная фильтрация — в основном на шаблоны в матрице рейтингов. Оба подхода могут использовать формальные предиктивные модели или эвристические алгоритмы, которые обычно ищут ближайших похожих пользователей или элементы. В дополнение к этим основным методам существует широкий спектр решений, объединяющий многоядерные алгоритмы в гибридные модели, расширяющий их для учета контекстных данных и вторичных целей оптимизации или дающий рекомендации в окружениях, где основные методы не являются оптимальными, например, из-за отсутствия данных персонализации. В следующих разделах мы детально проанализируем каждый из этих подходов.

Иерархия методов рекомендаций будет выглядеть иначе, если сосредоточиться на сценариях использования, а не на алгоритмических деталях и деталях реализации. Один из возможных вариантов классификации с этой позиции показан на рис. 5.5. Здесь все сценарии использования классифицируются в двух измерениях: по уровню персонализации — от неперсонализированных до сегментированных и полностью персонализированных, и степени использования контекстной информации — от контекстно-независимых до контекстно-зависимых. Основные методы рекомендаций, упоминавшиеся выше, фильтрация по содержимому и совместная фильтрация, в основном сосредоточены в персонализированном и контекстно-зависимом углу прямоугольника и генерируют рекомендации на основе

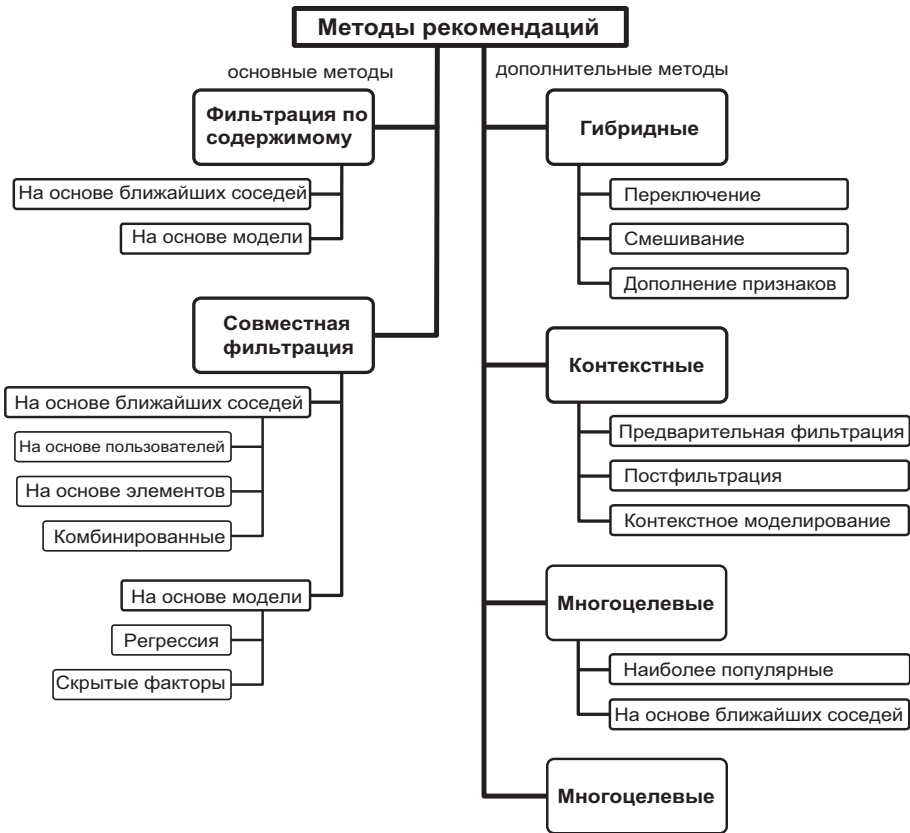


Рис. 5.4. Классификация методов рекомендаций с алгоритмической точки зрения

историй действий отдельных пользователей, то есть элементов, которые тот или иной пользователь оценил, просмотрел или купил в прошлом. В пользовательском интерфейсе подобные рекомендации часто отображаются в таких разделах, как *Вам могут понравиться*, *По вашей истории просмотров* или просто *Купить снова*. Рекомендации можно сделать еще более персонализированными, если учесть контекстную информацию, такую как местоположение пользователя, день недели или время суток. Примером этого класса методов рекомендаций может служить система рекомендаций ресторанов, предлагающая *Рестораны рядом с вами*, учитывающая историю действий пользователя и его местоположение. Альтернативный подход заключается в использовании неперсонализированных методов, основанных на глобальной статистике и свойствах элементов, а не на персональных профилях. В пользовательском интерфейсе эти рекомендации часто можно увидеть в таких разделах, как *Наиболее популярные*, *Пользуются спросом* или *Новые поступления*.

Обратите внимание, что персонализированные и неперсонализированные рекомендации можно смешивать. Например, персонализированные рекомендации, выбранные на основе истории действий пользователя, можно отсортировать по популярности или, как вариант, выбрать наиболее популярные элементы из категории продуктов, предпочитаемых пользователем. Наконец, неперсонализированные рекомендации можно уточнить по местоположению пользователя или атрибутам маркетингового канала. Например, страница с описанием продукта может включать разделы рекомендаций *Часто покупаются вместе* и *Похожие товары*, которые генерируются в контексте просматриваемого элемента.

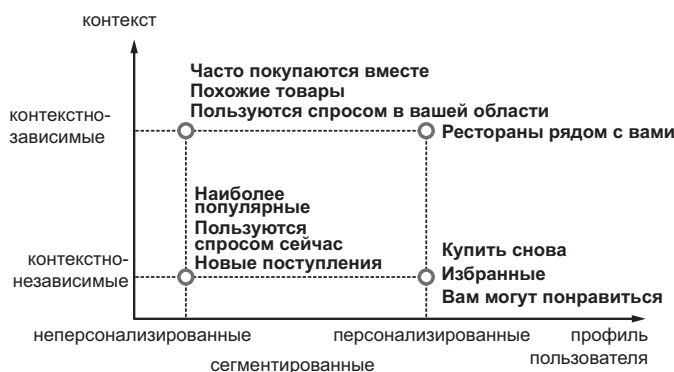


Рис. 5.5. Некоторые типичные сценарии использования и соответствующие методы рекомендаций

5.5. Фильтрация по содержимому

Первое семейство методов, которые мы рассмотрим, опирается в основном на каталог с данными (содержимое, или контент) и использует лишь малую часть информации, доступной в матрице рейтингов. По этой причине данная группа методов называется фильтрацией по содержимому. Основная идея довольно проста: берутся элементы, которые пользователь положительно оценил в прошлом, и рекомендуются элементы, похожие на них, как показано на рис. 5.6. Важное ограничение метода заключается в том, что мера сходства основана на содержимом элемента и не учитывает поведенческие данные, такие как сведения об элементах, которые часто приобретаются или высоко оцениваются другими пользователями.

¹ Подробное пояснение таких мер вы найдете в главе 4. Одним из наиболее типичных примеров является расстояние $TF \times IDF$ между текстовыми описаниями, как разъяснялось в разделе 4.3.5.

Фактически это означает, что такая рекомендательная система использует только одну строку матрицы рейтинга — профиль пользователя, которому адресованы рекомендации. Такое ограниченное использование информации о рейтингах часто компенсируется функцией сходства, использующей широкий спектр тщательно спроектированных признаков элементов. Затем сгенерированные рекомендации ранжируются в соответствии с их оценками сходства и, при необходимости, со значениями рейтингов соответствующих элементов в профиле. Например, если предположить, что элемент 1 имеет самый высокий рейтинг в примере, изображенном на рис. 5.6, то есть $r_{u1} > r_{u2}$ и $r_{u1} > r_{u3}$, элемент-кандидат, похожий на элемент 1, может получить более высокий ранг, чем элементы-кандидаты, одинаково похожие на элементы 2 или 3.

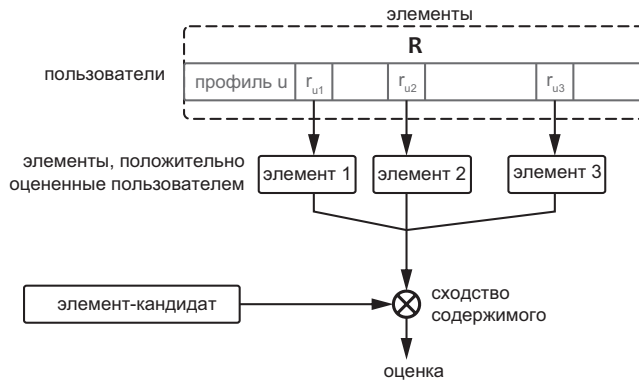


Рис. 5.6. Подход к фильтрации по содержанию на основе сходства

Однако метод фильтрации по содержанию, описанный выше, несколько ограничен, потому что фактически опирается на модель ближайших соседей — элемент-кандидат оценивается на основе среднего попарного сходства с элементами в профиле пользователя. Более общее и гибкое решение задачи фильтрации по содержанию заключается в использовании методов машинного обучения и оценке каждого элемента с помощью регрессионной или классификационной модели, обученной на профиле пользователя. Иначе говоря, для каждого пользователя создается специальная *модель профиля*, которая может предсказать, понравится ли пользователю данный элемент. Модель обучается на элементах из профиля: каждый элемент, оцененный пользователем, преобразуется в вектор признаков с помощью контент-анализатора, а соответствующий рейтинг используется в качестве обучающей метки. Затем каждый элемент-кандидат также преобразуется в вектор признаков и оценивается с помощью модели профиля, как показано на рис. 5.7. Наконец, путем ранжирования элементов-кандидатов в соответствии с прогнози-

руемыми рейтингами создается список рекомендаций и из него выбираются самые верхние. Обратите внимание, что подход к фильтрации по содержимому на основе сходства, изображенный на рис. 5.6, является частным случаем более общей схемы, представленной на рис. 5.7, полученной в предположении, что модель профиля использует классификатор по ближайшим соседям.

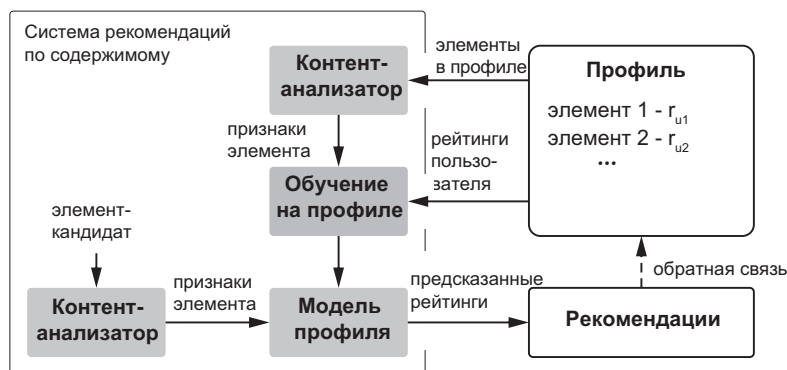


Рис. 5.7. Подход к фильтрации по содержимому на основе предсказания рейтинга

Фильтрация по содержимому имеет свои достоинства и недостатки. В числе основных достоинств систем рекомендаций по содержимому можно упомянуть:

- *Независимость от пользователя.* Фильтрация по содержимому использует только рейтинги, выставленные пользователем, для которого подготовлены рекомендации. Это может стать важным преимуществом, когда общее количество пользователей системы или собранных рейтингов невелико, как, например, в случае, когда новая рекомендательная система только вводится в эксплуатацию и располагает ограниченным объемом исторических данных. Эта проблема известна как *проблема холодного старта*. Вторым преимуществом независимости от пользователей является возможность выбора рекомендаций для пользователей с необычными вкусами, что трудно реализовать в системах, опирающихся на сходства пользователей и, как результат, склонным следовать средним вкусам.
- *Новые и редкие элементы.* Частным случаем проблемы холодного старта является рекомендация новых или редких элементов, имеющих небольшое количество рейтингов или не имеющих их вовсе. Алгоритм рекомендаций, опирающийся в основном на рейтинги, может не рекомендовать такие элементы, что негативно сказывается на охвате каталога. Фильтрация по содержимому не чувствительна к этой проблеме, поскольку основана на сходстве содержимого. Это особенно важно в контексте свойства длинного хвоста, которое мы

рассматривали выше, — каталог часто содержит много редких элементов, получивших небольшое количество рейтингов, даже в течение длительного периода времени. Та же проблема возникает в областях с быстро меняющимся ассортиментом, таких как магазины одежды, где бывает трудно накопить достаточное количество статистических данных о товарах.

- *Рекомендации по разным категориям.* Иногда бывает трудно или невозможно дать определенные рекомендации без учета содержимого. Например, трудно рекомендовать мебель людям, купавшим одежду, опираясь только на модель покупок, потому что число пользователей, которые покупают товары из обеих категорий одновременно, скорее всего, будет невелико [Ghani and Fano, 2002]. Как мы увидим далее, фильтрация по содержимому может оказаться подходящим решением для такого рода задач при надлежащем проектировании признаков.
- *Использование данных из каталога.* Фильтрация по содержимому позволяет использовать данные из каталога, которые являются важным источником информации для рекомендаций. Это контрастирует с некоторыми другими семействами алгоритмов рекомендаций, игнорирующими эти данные.
- *Интерпретируемость.* Рекомендации, сгенерированные системами на основе содержимого, часто легко интерпретируются или объясняются пользователями, поскольку оценки элементов основаны на определенных признаках содержимого. Например, рекомендательная система для фильмов может объяснить, что боевик рекомендован потому, что в прошлом пользователь высоко оценивал такие фильмы. Другие семейства методов рекомендаций могут давать результаты, которые гораздо труднее объяснить или интерпретировать.

С другой стороны, фильтрация по содержимому имеет ряд недостатков, которые часто являются обратной стороной ее преимуществ:

- *Проектирование признаков.* Очевидно, что проектирование признаков играет важную роль в фильтрации по содержимому, а качество прогнозирования рейтингов в значительной степени зависит от качества и полноты данных в каталоге и особенностей организации признаков. Проектирование признаков для элементов каталога — сложная задача даже для текстовых описаний продуктов из-за полисемии, синонимии и других проблем, которые мы рассматривали выше в контексте служб поиска. Задача может усложниться еще больше для изображений, фильмов или музыки, где нужно либо вручную снабдить каждый элемент тегами, такими как музыкальный жанр, или использовать продвинутое методы глубокого обучения. Управление содержимым и проектирование признаков действительно являются ключевыми проблемами в практических приложениях фильтрации по содержимому, и мы посвятим изучению этих аспектов много времени далее в этом разделе.

- *Новые пользователи.* Фильтрация по содержимому помогает решить проблему холодного старта для новых элементов, но она не в состоянии сгенерировать рекомендации для новых пользователей с пустыми профилями, что является второй разновидностью проблемы холодного старта.
- *Тривиальные рекомендации.* Одним из самых больших недостатков фильтрации по содержимому является склонность к тривиальным рекомендациям, то есть рекомендациям, которые не являются ни новыми, ни неожиданными. Это является прямым результатом оценки на основе содержимого, согласно которой предпочтение отдается тесно связанным элементам, таким как книги из одной серии.

Мы продолжим этот раздел более тщательным анализом методов рекомендаций, основанных на содержимом. Сначала рассмотрим два конкретных примера моделей профилей, а затем обсудим несколько продвинутых методов проектирования признаков, разработанных для рекомендательных систем в сфере розничной торговли.

5.5.1. Метод ближайших соседей

Как уже упоминалось, можно построить систему, генерирующую рекомендации по содержимому, используя в качестве модели профиля алгоритм k ближайших соседей (k Nearest Neighbor, k NN). Конкретизируем детали возможной реализации этого подхода. Прежде всего, обозначим набор элементов, оцененных пользователем u как I_u . Также можно предположить, что каждый элемент j представлен документом d_j с одним или несколькими атрибутами или полями, аналогично ситуации со службами поиска. Следовательно, совокупность элементов I_u соответствует совокупности документов D_u . Для каждого элемента-кандидата i можно вычислить метрики сходства между представлением документа d_i и каждым из документов в D_u . Обозначим k документов в D_u , имеющих наибольшее сходство с d_i , как

$$\{d_1^{ui}, \dots, d_k^{ui}\} \subset D_u. \quad (5.14)$$

Эти документы являются k ближайшими соседями d_i в соответствии с некоторой мерой сходства. Рейтинг пользователя для элемента i можно предсказать как средний рейтинг ближайших соседей:

$$\hat{r}_{ui} = \frac{1}{k} \sum_{t=1}^k r_u(d_t^{ui}), \quad (5.15)$$

где $r_u(d)$ — рейтинг элемента, соответствующего документу d . Эту оценку можно уточнить взвешиванием рейтингов оценками сходства:

$$\hat{r}_{ui} = \frac{1}{k} \sum_{t=1}^k r_u(d_t^{ui}) \cdot \text{sim}(d_i, d_t^{ui}). \quad (5.16)$$

Мера сходства обычно вычисляется с использованием приемов, разработанных для служб поиска. Один из самых популярных — использование базового векторного пространства модели: к текстовым полям документов применяется стемминг, из них удаляются стоп-слова, после этого каждое поле документа преобразуется в вектор термов, затем между соответствующими полями вычисляется расстояние в соответствии с моделью $TF \times IDF$ и, наконец, с помощью некоторой функции смешивания сигналов оценки для разных полей объединяются в окончательную оценку сходства. Это в точности тот же процесс оценки документов, что используется в службах поиска; единственное отличие — оценка сходства $TF \times IDF$ вычисляется для пары документов, а не для документа и запроса.

Второй популярный вариант — использование модели скрытой темы вместо модели базового векторного пространства: с помощью латентно-семантического анализа (LSA) или латентного распределения Дирихле (LDA) каждое поле документа представляется как вектор в пространстве скрытых тем, между соответствующими полями вычисляется расстояние как косинусное расстояние между соответствующими векторами и затем оценки для отдельных полей объединяются в итоговую оценку. И снова этот подход почти идентичен методам поиска LSA и LDA, которые мы рассматривали выше. Принято считать, что метод LDA превосходит LSA и базовые модели $TF \times IDF$ [Falk, 2017]. Он с успехом использовался в некоторых крупных промышленных системах, таких как механизм рекомендации статей в *The New York Times* [Spangher, 2015]. С другой стороны, согласно некоторым отчетам, LSA может превзойти LDA в определенных приложениях, таких как рекомендательные системы для фильмов [Bergamaschi et al., 2014]. Очевидно, что результаты сильно зависят от используемых наборов данных и методологии оценки качества.

5.5.2. Наивный байесовский классификатор

Второй подход к фильтрации по содержанию, который мы рассмотрим, был разработан для рекомендаций книг [Mooney and Roy, 1999]. В отличие от регрессии по ближайшим соседям, в качестве модели профиля для прогнозирования рейтингов этот метод использует не эвристические метрики сходства, а наивный байесовский классификатор — стандартный алгоритм классификации текста.

Прежде всего предположим, что каждый элемент каталога имеет несколько текстовых атрибутов. Например, книга может иметь такие атрибуты, как название, авторы, обзор, опубликованные отзывы, комментарии клиентов, похожие книги и похожие авторы. Применим к атрибутам лексемизацию и стемминг, удалим стоп-слова, а затем смоделируем каждый атрибут как мешок слов, то есть вектор, каждый элемент которого соответствует слову со значением, равным количеству

вхождений этого слова в текст атрибута. То есть каждый элемент представлен документом с несколькими полями и каждое поле является моделью мешка слов соответствующего атрибута.

Далее мы должны создать модель профиля. Напомню, что конечной целью рекомендаций по содержанию является ранжирование элементов каталога для каждого пользователя в порядке его предпочтений. Эту задачу можно решить, построив бинарный классификатор, оценивающий две вероятности: вероятность положительной и вероятность отрицательной оценки элемента пользователем. Отношение между этими двумя вероятностями указывает, будет ли элемент оценен положительно, и, следовательно, его можно использовать в качестве оценки для ранжирования рекомендуемых элементов. Предположим, что пользователь оценивает элементы по дискретной шкале от 1 до r_{max} и все оценки ниже $r_{max}/2$ интерпретируются как отрицательные, а оценки выше $r_{max}/2$ — как положительные. Например, для шкалы рейтингов от 1 до 10 звезд рейтинги 1–5 считаются отрицательными (не нравится), а рейтинги 6–10 — положительными (нравится).

Напомню, что основная идея наивного байесовского текстового классификатора заключается в оценке вероятности принадлежности документа d к некоторому классу c с использованием условных вероятностей вхождения слов документа w в документ класса c , в предположении, что эти условные вероятности независимы. Этот подход можно выразить с помощью следующей формулы:

$$\Pr(c_j | d) = \frac{\Pr(c_j)}{\Pr(d)} \prod_{w_i \in d} \Pr(w_i | c_j), \quad (5.17)$$

где c_j — класс документа, который в нашем случае является либо отрицательным c_0 , либо положительным c_1 , $\Pr(c_j)$ — эмпирическая вероятность класса c_j в обучающих данных (доля документов, принадлежащих классу), и $\Pr(w_i | c_j)$ эмпирическая условная вероятность слова w_i (доля документов класса c_j , содержащих это слово). Это основное правило Байеса требуется распространить на случай нескольких полей, имеющих в каждом документе с описанием элемента. Предположив, что каждый документ имеет F полей и каждое поле f_{qm} является текстовым фрагментом, содержащим $|f_{qm}|$, мы можем переписать формулу 5.17 для апостериорной вероятности класса следующим образом:

$$\Pr(c_j, d) = \frac{\Pr(c_j)}{\Pr(d)} \prod_{m=1}^F \prod_{w_i \in f_m} \Pr(w_i | c_j, f_m), \quad (5.18)$$

Оценку ранжирования элемента тогда можно определить как

$$\text{score}(d) = \frac{\Pr(c_1 | d)}{\Pr(c_0 | d)}, \quad (5.19)$$

и, соответственно, отсортировать элементы в списке рекомендаций в порядке ее убывания.

Наш следующий шаг — оценка вероятностей в формуле 5.18 на основе профиля пользователя, то есть на основе элементов, оцененных пользователем. Как говорилось выше, рейтинги пользователей устанавливаются по шкале от 1 до r_{\max} . В случае, если пользователь оценил Q элементов, мы можем отобразить каждый рейтинг в две вспомогательные переменные, для положительного и отрицательного классов соответственно:

$$\alpha_{q1} = \frac{r_q - 1}{r_{\max} - 1}, \quad q = 1, \dots, Q, \quad (5.20)$$

$$\alpha_{q0} = 1 - \alpha_{q1}, \quad q = 1, \dots, Q, \quad (5.21)$$

где r_q — исходный рейтинг в профиле пользователя. Обратите внимание, что во всех уравнениях последовательно опускается индекс пользователя u , потому что алгоритм использует только профиль активного пользователя. Вероятность класса можно оценить следующим образом:

$$\Pr(c_j) = \frac{1}{Q} \sum_{q=1}^Q \alpha_{qj}, \quad j = 0, 1. \quad (5.22)$$

Условные вероятности слов должны оцениваться для каждого поля документа отдельно. Если число вхождений слова w_i в поле m документа q обозначить как $n_{qm}(w_i)$, тогда условную вероятность этого слова можно оценить как

$$\Pr(w_i | c_j, \text{field} = m) = \sum_{q=1}^Q \alpha_{qj} \cdot \frac{n_{qm}(w_i)}{L_{jm}}, \quad m = 1, \dots, F, \quad (5.23)$$

где L_{jm} — общая взвешенная длина текстов в поле m для класса j :

$$L_{jm} = \sum_{q=1}^Q \alpha_{qj} \cdot |f_{qm}|, \quad m = 1, \dots, F. \quad (5.24)$$

Длина поля определяется как количество слов в его представлении как мешка слов. Эти оценки позволяют оценить апостериорные вероятности класса документа из уравнения 5.18 и в конечном итоге оценить сами элементы. Обратите внимание, что вероятности $\Pr(d)$ можно игнорировать, поскольку они нейтрализуют друг друга в формуле оценки 5.19.

ПРИМЕР 5.1

Рассмотрим конкретный пример, чтобы лучше понять, как наивный байесовский классификатор может давать рекомендации и какие ограничения имеет этот подход. Рассмотрим книжный онлайн-магазин, в каталоге которого каждая книга представлена документом с двумя полями: *title* (название) и *synopsis* (описание). Допустим у нас имеется профиль пользователя, оценившего две книги по шкале от 1 до 10, и создадим для него модель профиля. Исходный профиль выглядит следующим образом:

Book 1

Title: Machine learning for predictive data analytics

Synopsis: Detailed treatment of data analytics applications
including price prediction and customer behavior

Rating: 8

Book 2

Title: Machine learning for healthcare and life science

Synopsis: Case studies specific to the challenges of
working with healthcare data

Rating: 3

Преобразовав поля в мешки слов и удалив стоп-слова, получим:

```
title1 : (machine, learning, predictive, data, analytics)
synopsis1 : (detailed, treatment, data, analytics, applications,
            including, price, prediction, customer, behavior)
title2 : (machine, learning, healthcare, life, science)
synopsis2 : (case, studies, specific, challenges, working,
            healthcare, data)
```

Вычислим значения близости классов согласно формулам 5.20 и 5.21:

$$\alpha_{11} = \frac{8-1}{9} = \frac{7}{9}, \quad (5.25)$$

$$\alpha_{10} = 1 - \alpha_{11} = \frac{2}{9}, \quad (5.26)$$

$$\alpha_{21} = \frac{3-1}{9} = \frac{2}{9}, \quad (5.27)$$

$$\alpha_{20} = 1 - \alpha_{21} = \frac{7}{9}. \quad (5.28)$$

Используем эти значения для оценки вероятностей классов в соответствии с выражением 5.22. Так как первая книга понравилась пользователю (рей-

тинг 8 из 10), а вторая не понравилась (рейтинг 3 из 10), вероятности будут равны:

$$\begin{aligned}\Pr(c_0) &= \frac{1}{2}(\alpha_{10} + \alpha_{20}) = \frac{1}{2}, \\ \Pr(c_1) &= \frac{1}{2}(\alpha_{11} + \alpha_{21}) = \frac{1}{2}.\end{aligned}\tag{5.29}$$

Вычислив взвешенные длины полей по формуле 5.24, получим:

$$\begin{aligned}L_{0, \text{title}} &= \alpha_{10}|\text{title}_1| + \alpha_{20}|\text{title}_2| = \frac{2}{9} \cdot 5 + \frac{7}{9} \cdot 5 = 5, \\ L_{1, \text{title}} &= \alpha_{11}|\text{title}_1| + \alpha_{21}|\text{title}_2| = 5, \\ L_{0, \text{synopsis}} &= \alpha_{10}|\text{synopsis}_1| + \alpha_{20}|\text{synopsis}_2| = \frac{23}{3}, \\ L_{1, \text{synopsis}} &= \alpha_{11}|\text{synopsis}_1| + \alpha_{21}|\text{synopsis}_2| = \frac{28}{3}.\end{aligned}\tag{5.30}$$

Наконец, можно оценить условные вероятности слов, используя выражение 5.23. В качестве иллюстрации оценим условную вероятность слова *price* в поле *synopsis* для отрицательного класса:

$$\begin{aligned}\Pr(\text{price} | c = 0, \text{field} = \text{synopsis}) &= \\ &= \alpha_{10} \frac{n_{1, \text{synopsis}}(\text{price})}{L_{0, \text{synopsis}}} + \alpha_{20} \cdot \frac{n_{2, \text{synopsis}}(\text{price})}{L_{0, \text{synopsis}}} = \\ &= \frac{2}{9} \cdot \frac{1}{23/3} + \frac{7}{9} \cdot \frac{0}{23/3} = \frac{2}{69}.\end{aligned}\tag{5.31}$$

Оценив вероятности для всех сочетаний слов, классов и полей, получим результаты, приведенные в табл. 5.1. Эта таблица фактически является моделью профиля, которая будет вычисляться системой рекомендаций заранее, храниться и использоваться для ранжирования рекомендаций в соответствии с формулами 5.18 и 5.19.

Таблица 5.1 содержит некоторые полезные сведения о логике рекомендательной системы на основе наивного байесовского классификатора. Первое, что бросается в глаза, — слова, встречающиеся в положительно и отрицательно оцененных элементах, нейтрализуют друг друга. Например, в названиях обеих книг есть слова *machine learning* (машинное обучение), поэтому каждое из них имеет равное значение вероятности для положительных и отрицательных

Таблица 5.1. Пример модели профиля, полученной с применением наивного байесовского классификатора

	title		synopsis	
	c = 0	c = 1	c = 0	c = 1
ANALYTICS	0,044	0,160	0,029	0,083
APPLICATIONS	0,000	0,000	0,029	0,083
BEHAVIOR	0,000	0,000	► 0,029	► 0,083
CASE	0,000	0,000	0,100	0,024
CHALLENGES	0,000	0,000	0,100	0,024
CUSTOMER	0,000	0,000	0,029	0,083
DATA	0,044	0,160	0,130	0,110
DETAILED	0,000	0,000	0,029	0,083
HEALTHCARE	► 0,160	► 0,044	0,100	0,024
INCLUDING	0,000	0,000	0,029	0,083
LEARNING	► 0,200	► 0,200	0,000	0,000
LIFE	0,160	0,044	0,000	0,000
MACHINE	► 0,200	► 0,200	0,000	0,000
PREDICTION	0,000	0,000	0,029	0,083
PREDICTIVE	0,044	0,160	0,000	0,000
PRICE	0,000	0,000	0,029	0,083
SCIENCE	0,160	0,044	0,000	0,000
SPECIFIC	0,000	0,000	0,100	0,024
STUDIES	0,000	0,000	0,100	0,024
TREATMENT	0,000	0,000	0,029	0,083
WORKING	0,000	0,000	0,100	0,024

классов, и эти значения будут нейтрализовать друг друга в соотношении 5.19, используемом для оценки. Во-вторых, слова, присутствующие в атрибутах отрицательно оцененной книги (например, *healthcare* (здравоохранение)), интерпретируются как негативные сигналы, в том смысле что их вероятности для отрицательного класса выше, чем для положительного. Аналогично некоторые другие слова (например, *behavior* (поведение)) интерпретируются

как положительные сигналы. На практике такое толкование не всегда бывает правильным. В нашем примере пользователю не понравилась вторая книга о машинном обучении для здравоохранения. Действительная причина нам не известна — может быть, эта конкретная книга плохо написана или сфера здравоохранения нерелевантна для пользователя. Если предположить, что пользователь оценил книгу после покупки и прочтения, первое объяснение выглядит более вероятным, потому что пользователь почти наверняка знал, что книга относится к сфере здравоохранения и выбрал ее сознательно. Наивная байесовская модель, однако, интерпретирует слово *healthcare* как отрицательный сигнал, поэтому все книги с этим словом в названии получают пониженную оценку. Эта ограниченная способность различать качество и релевантность содержимого является одним из основных недостатков фильтрации по содержанию. Совместная фильтрация, как мы увидим далее, использует другой подход к задаче и уделяет больше внимания сигналам качества элементов.

5.5.3. Проектирование признаков для фильтрации по содержанию

Основная идея фильтрации по содержанию заключается в создании регрессионной или классификационной модели, оценивающей содержание элемента. Очевидно, что такой подход требует тщательного проектирования признаков, поскольку качество классификации в значительной степени зависит от доступных атрибутов изделия и качества их моделирования. Тривиальные атрибуты часто приводят к тривиальным или бессмысленным рекомендациям, тогда как тщательно спроектированные признаки способны помочь рекомендательной системе точно предсказать решения, выбираемые пользователями. Проиллюстрируем эту проблему на примере из области торговли одеждой [Ghani and Fano, 2002; Ghani et al., 2006]. Рассмотрим пользователя, который приобрел и оценил несколько предметов одежды, таких как платья, блузки или пальто. Можно ожидать, что типичная система управления информацией о продуктах будет способна представить некоторую базовую информацию о каждом из этих элементов, например категорию, цену и цвет. Простая система рекомендаций, которая вычисляет сходства между элементами с помощью подобных атрибутов, скорее всего, будет рекомендовать элементы из тех же категорий, в том же ценовом диапазоне и того же цвета. Такой подход не обязательно дает плохие рекомендации, но он имеет, по крайней мере, два недостатка. Во-первых, выбор клиента в значительной степени зависит от личности, вкуса и образа жизни. Клиенты склонны думать об одежде с точки зрения стиля и функциональности, выбирая между свободным и формальным стилем, спортивным и деловым, консервативным и ярким. И пользователей,

и одежду можно описать с точки зрения таких *психографических особенностей*, и рекомендации можно сгенерировать на основе близости пользователя к определенным стилям и вкусам. Рекомендательная система, использующая только базовые атрибуты, такие как категория продукта и цена, обычно не распознает эти скрытые сходства. Во-вторых, некоторые виды рекомендаций, например рекомендации продуктов из разных категорий, в принципе трудно сгенерировать, если доступны только простые признаки. Одна из причин заключается в том, что разбивка на атрибуты для продуктов в разных категориях может различаться, поэтому иногда трудно определить метрику сходства модели профиля, с помощью которой можно сравнивать и оценивать элементы в разных категориях. Например, большой универсам может продавать одежду, кухонную утварь и мебель. Однако очень сложно рекомендовать мебель, опираясь на покупки в отделе одежды, потому что одежда и мебель имеют разные атрибуты с разной семантикой. Например, атрибут *размер* имеет совершенно разное значение для платья и для кровати. И снова решить проблему могут помочь психографические особенности, потому что пользователям, покупающим консервативную одежду, можно рекомендовать мебель консервативного стиля и т. д.

Одно из наблюдений, которые можно сделать, заключается в том, что текстовые атрибуты продукта, такие как название, описание и отзывы клиентов, часто несут неявный сигнал о психографических характеристиках продукта. Мерчендайзеры, придумывающие названия продуктам и составляющие их описания, часто сознательно выбирают определенные слова, такие как *стильный*, *сексуальный* или *роскошный*, чтобы повысить привлекательность продукта для определенной аудитории. Этот факт можно использовать для оценки близости продукта к определенным психографическим характеристикам и определения соответствующих признаков. Далее эти признаки можно использовать для обучения и использования модели профиля. Более конкретно для извлечения неявных психографических характеристик можно использовать следующий метод [Ghani and Fano, 2002]:

- Сначала определяется набор признаков продукта, которые извлекаются с использованием знаний о предметной области. Примеры таких характеристик для одежды приводятся в табл. 5.2.
- Затем экспертами в предметной области подмножеству элементов вручную присваиваются метки в соответствии с признаками, определенными на предыдущем шаге. Это подмножество используется для обучения моделей классификации, предсказывающих психографические метки на основе текстовых описаний и других стандартных атрибутов продуктов, таких как название бренда и размер. Например, для идентификации слов в описаниях продуктов, указывающих на высокую степень формальности или современности, можно использовать наивный байесовский классификатор.

Таблица 5.2. Примеры психографических характеристик для одежды

Характеристика	Значения характеристики
Возрастная группа	Наиболее соответствующая возрастная группа: дети, подростки, взрослые и т. д.
Назначение	Типичный сценарий использования продукта: вечерняя одежда, спортивная, деловая повседневная, деловая формальная и т. д.
Степень формальности	От неформальной до очень формальной
Степень консерватизма	От очень консервативной одежды, такой как серый костюм, до крикливой и вычурной
Степень спортивности	От небрежной или формальной до исключительно спортивной
Степень современности	От неувядающей классики до последних модных веяний
Степень известности бренда	От неизвестного или непопулярного до очень известного бренда

- Обученные модели классификации используются для присваивания меток остальным элементам. Это позволяет мерчендайзерам с минимальными усилиями снабдить метками даже очень большие и часто изменяющиеся каталоги.

Стоит отметить, что некоторые рекомендательные системы полностью полагаются на атрибуты элементов, созданные вручную. Одним из известных примеров является Pandora Internet Radio, служба потоковой передачи музыки и рекомендаций. В Pandora пользуются услугами профессиональных музыкальных аналитиков, чтобы вручную снабдить каждую песню в своем каталоге 450 признаками, такими как *Детская* или *Детский вокал* и *Мелодическая артикуляция от ясной до неразборчивой*. Этот анализ, известный как Music Genome Project, требует значительных усилий, потому что каталог содержит сотни тысяч песен, а классификация одной песни занимает около 20 минут [Walker, 2009]. Однако эти метаданные являются основным активом Pandora и основным конкурентным преимуществом на рынке услуг по поиску музыки.

5.6. Введение в совместную фильтрацию

Фильтрация по содержанию пытается аппроксимировать вкусы и суждения пользователей с помощью меры сходства содержимого элементов каталога. Основным недостатком этого подхода связан с отсутствием простого способа выражения вкусов

человека в виде простых атрибутов, поэтому для достижения хороших результатов часто требуется ручная маркировка продуктов и продвинутое проектирование признаков. С другой стороны, матрица рейтингов содержит много информации о вкусах и суждениях пользователей. Действительно, каждый известный рейтинг можно интерпретировать как атрибут продукта, заданный вручную, и, следовательно, сбор рейтингов и другой обратной связи пользователей можно рассматривать как коллективный подход к маркировке продуктов психографическими характеристиками. Фильтрация по содержанию не полностью использует эту ценную информацию, потому что рекомендации создаются с помощью единственной модели профиля. Рассуждая в этом направлении, мы неизбежно подходим к другому семейству методов рекомендаций, известному как совместная фильтрация.

Термин *совместная фильтрация* был придуман разработчиками Tapestry, рекомендательной системы для новостей и статей, созданной в Xerox PARC в 1992 году [Goldberg et al., 1992; Terry, 1993]. В контексте Tapestry под совместной фильтрацией подразумевалось, что пользователи могли посылать отзывы в ответ на электронные письма с новостями и выражать предпочтение сообщениям, исходя из отзывов других пользователей. По сути это была функциональная возможность фильтрации электронной почты, а не алгоритм рекомендаций. Однако идея выражения предпочтений в отношении рекомендаций по отзывам других пользователей приобрела большую популярность и привела к разработке новых методов рекомендаций, которые использовали этот подход и его методы в промышленных рекомендательных системах, включая разработанные в компаниях Amazon и Netflix. Значение термина «совместная фильтрация» тоже изменилось и теперь в большей степени относится к прогнозированию рейтингов на основе информации, доступной в матрице рейтингов. Совместная фильтрация в этом новом, более узком смысле является чистой задачей восстановления матрицы. Следовательно, методы совместной фильтрации, по сути, являются алгоритмами восстановления матрицы, которые используют матрицу рейтингов в качестве единственного входа. За кулисами совместная фильтрация использует предиктивную модель, которая по взаимодействиям между пользователями и элементами, известными из матрицы рейтингов, предсказывает рейтинг для данной пары пользователь/элемент на основе рейтингов, выставленных в прошлом похожими пользователями похожим элементам.

Главное преимущество методов совместной фильтрации заключается в их способности давать рекомендации, основанные только на шаблонах и сходствах, доступных в матрице рейтингов, без какой-либо дополнительной информации об элементах в каталоге. Это делает методы совместной фильтрации гораздо более универсальными, чем фильтрация по содержанию, которая требует специальных знаний, данных и приложения усилий по проектированию признаков. Что еще более важно, совместная фильтрация неявно учитывает психографические профили пользователей и элементов, потому что рейтинги отражают человеческие вкусы

и суждения. Это помогает вырабатывать нетривиальные рекомендации. С другой стороны, совместная фильтрация имеет ряд недостатков:

- *Разреженность рейтингов.* Рекомендательная система на основе совместной фильтрации требует достаточного количества известных и заслуживающих доверия оценок. Если матрица рейтингов слишком разрежена, построение надежной модели прогнозирования рейтингов может оказаться затрудненным или невозможным.
- *Новые пользователи и элементы.* Совместная фильтрация прогнозирует рейтинги для данного пользователя или элемента на основе известных рейтингов для этого пользователя или элемента. То есть совместная фильтрация плохо подходит для новых пользователей и элементов или для пользователей и элементов с очень небольшим количеством известных рейтингов. То есть совместная фильтрация более подвержена проблеме холодного старта, чем фильтрация по содержанию, которая кроме рейтингов использует информацию из содержимого.
- *Предвзятость в пользу популярности.* Совместная фильтрация дает рекомендации, основанные на типичных шаблонах в матрице рейтингов, поэтому она изначально склонна к выбору популярных элементов и стандартных вариантов. Это ограничивает возможность получения нетривиальных рекомендаций и рекомендаций для пользователей с необычными вкусами.
- *Стандартизация продуктов.* Совместная фильтрация в принципе способна распознать элементы, часто приобретаемые вместе, но вообще она рассматривает каждый элемент как непрозрачную и независимую сущность. Это может создать определенные проблемы для продуктов со сложной внутренней структурой, таких как предметы одежды разных размеров, продукты, подгоняемые под требования клиентов, или продукты, обновляющиеся с течением времени.
- *Знание предметной области.* Как уже отмечалось, одним из основных преимуществ совместной фильтрации является возможность работы с абстрактной матрицей рейтингов без каких-либо предположений о природе элементов и их атрибутов. В целом это верно, но иногда методы совместной фильтрации могут потребовать сделать определенные предположения, относящиеся к конкретной области. Например, рекомендательная система на основе совместной фильтрации может предполагать или не предполагать изменение вкусов клиентов с течением времени и, следовательно, может учитывать или не учитывать давность рейтинга.

Алгоритмы совместной фильтрации обычно делятся на две подгруппы: методы на основе близости и методы на основе модели. Методы, основанные на близости (также известные как анамнестические (memory-based) методы), предсказывают

неизвестные рейтинги для данного пользователя или элемента с помощью метода ближайшего соседа, то есть путем определения наиболее похожих пользователей или элементов и усреднения известных рейтингов из их записей. Методы, основанные на моделях, выходят за рамки подхода ближайших соседей и используют другие, обычно более сложные предиктивные модели. Хотя алгоритм ближайших соседей тоже можно рассматривать как своеобразную предиктивную модель (что делает границу между двумя категориями немного размытой), имеет смысл разделить их вследствие большой практической ценности методов на основе близости. Мы подробно рассмотрим оба подхода в последующих разделах.

5.6.1. Базовые оценки

Большинство моделей совместной фильтрации, используемых на практике, способны фиксировать относительно сложные взаимодействия между пользователями и элементами путем распознавания сложных шаблонов в матрице рейтингов. Прежде чем приступить к изучению этих моделей, важно отметить, что наблюдаемые рейтинги обычно следуют нескольким простым, но мощным шаблонам, которые можно обнаруживать с помощью относительно простой модели. Эта простая модель способна генерировать базовые оценки рейтингов, которые затем могут использоваться в качестве строительных блоков для создания более продвинутых методов совместной фильтрации.

Типичная матрица рейтингов демонстрирует выраженную *предвзятость* в отношении пользователей и элементов — некоторые пользователи систематически дают более высокие (или низкие) рейтинги, чем другие, и некоторые элементы систематически получают более высокий рейтинг, чем другие [Koren, 2009; Ekstrand et al., 2011]. Это можно объяснить тем, что некоторые пользователи более или менее критичны, чем другие, и элементы, конечно, отличаются своими качествами. Мы можем учесть эти систематические эффекты, определив базовую оценку неизвестного рейтинга пользователя r_{ui} как

$$b_{ui} = \mu + b_u + b_i, \quad (5.32)$$

где μ — общий средний рейтинг в матрице рейтингов R , b_u — наблюдаемое отклонение пользователя u от среднего, а b_i — наблюдаемое отклонение элемента i от среднего. На практике смещенность в оценках оказывает сильное влияние и, следовательно, базовая оценка, определяемая уравнением 5.32, обладает значительной предсказательной способностью. Хотя эта модель учитывает только усредненные эффекты, она может помочь уравновесить предвзятость и изолировать сигнал, представляющий взаимодействия между пользователем и элементом, который может быть воспринят более специализированными моделями.

Смещения μ , b_u и b_i можно оценить друг за другом как средние остаточные ошибки предыдущей оценки. Это означает, что сначала вычисляется μ , и только потом оцениваются смещения элементов b_i :

$$b_i = \frac{1}{|U_i|} \sum_{i \in U_i} (r_{ui} - \mu), \quad (5.33)$$

где U_i — множество пользователей, оценивших элемент i . Затем оценивается смещенность (предвзятость) пользователя:

$$b_u = \frac{1}{|I_u|} \sum_{u \in I_u} (r_{ui} - \mu - b_i), \quad (5.34)$$

где I_u — элементы, оцененные пользователем u . Оценки, вычисленные по формулам 5.33 и 5.34, могут быть неустойчивыми в случае разреженной матрицы рейтингов, когда для пользователя или элемента доступно только небольшое количество рейтингов. Стабильность оценок можно улучшить добавлением параметров регуляризации λ_1 и λ_2 :

$$\begin{aligned} b_i &= \frac{1}{|U_i| + \lambda_1} \sum_{i \in U_i} (r_{ui} - \mu), \\ b_u &= \frac{1}{|I_u| + \lambda_2} \sum_{u \in I_u} (r_{ui} - \mu - b_i). \end{aligned} \quad (5.35)$$

Параметры регуляризации уменьшают смещения b_i и b_u , когда пользователь или элемент имеет небольшое количество рейтингов, поэтому базовая оценка, описанная уравнением 5.32, становится ближе к глобальному среднему и меньше зависит от ненадежных оценок смещения.

Значения смещенности можно оценить точнее, решив следующую задачу наименьших квадратов [Koren, 2009]:

$$\min_{b_i, b_u} \sum_{i, u \in R} (r_{ui} - \mu - b_i - b_u) + \lambda \cdot \left(\sum_u b_u^2 + \sum_i b_i^2 \right), \quad (5.36)$$

где R — обучающий набор известных рейтингов и λ — параметр регуляризации. Это простая задача оптимизации, которая решается с помощью стандартных методов, таких как стохастический градиентный спуск. Преимущество этого подхода заключается в том, что выражение 5.36 можно легко изменить и расширить, включив в него дополнительные ограничения и переменные, например, описывающие эффекты, проявляющиеся с течением времени.

ПРИМЕР 5.2

Для иллюстрации вычисления базовых оценок и других методов совместной фильтрации нам понадобится образец матрицы рейтингов. Я предлагаю использовать пример с рейтингами фильмов, ставший очень популярным после конкурса Netflix Prize. Кстати отмечу, что все методы совместной фильтрации, представленные в этой главе, являются универсальными — они не зависят от природы элементов, поэтому названия фильмов выбраны исключительно из-за удобства чтения и их можно заменить любыми другими продуктами, такими как бакалея или одежда.

Примером нам послужит матрица рейтингов с шестью фильмами и шестью пользователями, представленная в табл. 5.3. Рейтинги присвоены по 5-звездочной шкале, где 1 — самый низкий возможный рейтинг и 5 — самый высокий. Матрица включает 28 известных и 8 отсутствующих рейтингов, то есть она очень плотная в сравнении с матрицами рейтингов в реальной жизни, где может отсутствовать до 99 % возможных рейтингов. В представленной матрице легко заметить несколько шаблонов. Во-первых, первым трем пользователям определенно больше нравятся драматические фильмы (*Forrest Gump* («Форрест Гамп»), *Titanic* («Титаник») и *The Godfather* («Крестный отец»)), чем боевики (*Batman* («Бэтмен»), *The Matrix* («Матрица») и *Alien* («Чужой»)). Последние три пользователя, напротив, по всей видимости, предпочитают боевики. Далее можно отметить, что пользователь 3 щедро присваивает высокие рейтинги большинству фильмов, тогда как пользователь 2 кажется более критичным. Будем надеяться, что хорошая модель совместной фильтрации распознает эти шаблоны и сделает адекватные прогнозы.

Таблица 5.3. Пример матрицы рейтингов для службы рекомендаций фильмов

	Forrest Gump	Titanic	The Godfather	Batman	The Matrix	Alien
Пользователь 1	5	4	—	1	2	1
Пользователь 2	4	—	3	1	1	2
Пользователь 3	—	5	5	—	3	3
Пользователь 4	2	—	1	4	5	4
Пользователь 5	2	2	2	—	4	—
Пользователь 6	1	2	1	—	5	4

Теперь вычислим базовые оценки для отсутствующих рейтингов. Определяя глобальное среднее и значения смещений по формулам 5.33 и 5.34, получаем:

$$\begin{aligned}\mu &= 2,82, \\ b_i &= (-0,02 + 0,42 - 0,42 - 0,82 + 0,51 - 0,02), \\ b_u &= (-0,23 - 0,46 + 1,05 + 0,53 - 0,44 - 0,31).\end{aligned}\tag{5.37}$$

Обратите внимание: несмотря на простоту, эти коэффициенты правильно отражают, что пользователь 3 склонен давать высокие рейтинги (смещение +1,05), а фильм *Batman* в целом оценивается довольно низко (смещение -0,82). Подставив этот результат в формулу базовой оценки 5.32, получим предиктивные значения отсутствующих рейтингов, показанные в табл. 5.4. Обратите внимание, что в целом результат не соответствует нашим ожиданиям о сходстве между пользователями и жанрами фильмов.

Таблица 5.4. Базовые оценки рейтингов

	Forrest Gump	Titanic	The Godfather	Batman	The Matrix	Alien
Пользователь 1	5	4	[2,16]	1	2	1
Пользователь 2	4	[2,78]	3	1	1	2
Пользователь 3	[3,85]	5	5	[3.05]	3	3
Пользователь 4	2	[3,78]	1	4	5	4
Пользователь 5	2	2	2	[1,55]	4	[2,35]
Пользователь 6	1	2	1	[1,68]	5	4

5.7. Совместная фильтрация на основе близости

Совместная фильтрация на основе близости опирается на меру сходства между пользователями или элементами по рейтингам, которые имеют два пользователя или два элемента. Эти два случая — сходство по пользователям и сходство по элементам — различны, но имеют много общего.

Рассмотрим сначала подход к фильтрации по близости пользователей и изображенный на рис. 5.8. Напомню, что цель рекомендательной системы состоит

в предсказании рейтингов, которые определенный пользователь дал бы разным элементам каталога, а затем — в создании списка рекомендаций, выбрав и оценив элементы с самыми высокими прогнозируемыми рейтингами. Если предположить, что пользователь уже оценил некоторые элементы в каталоге и соответствующая строка в матрице рейтингов содержит некоторые известные значения, можно попробовать найти еще пользователей, давших тем же элементам оценки с той же эмоциональной окраской, то есть которым понравились элементы, положительно оцененные данным пользователем, и не понравились отрицательно оцененные. Главная идея совместной фильтрации на основе близости заключается в том, что такие пользователи, скорее всего, будут иметь те же вкусы и предпочтения, что и данный пользователь, поэтому их рейтинги, выставленные в прошлом, можно использовать для прогнозирования будущих рейтингов данного пользователя. Следовательно, система сможет рекомендовать элементы, не оцененные данным пользователем, но положительно оцененные, по крайней мере, некоторыми близкими пользователями. Прогнозируемые оценки для этих элементов можно получить усреднением рейтингов, присвоенных близкими пользователями.

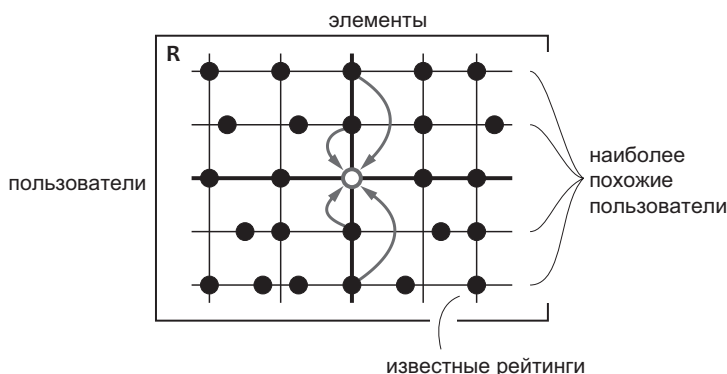


Рис. 5.8. Совместная фильтрация по близости пользователей

Подход к фильтрации по близости элементов, изображенный на рис. 5.9, структурно похож на только что описанный, с той лишь разницей, что пользователи (строки) заменяются элементами (столбцами). Чтобы спрогнозировать рейтинги для данного элемента, сначала нужно найти элементы, похожие на данный, то есть одинаково оцениваемые одними и теми же пользователями. Затем рейтинг, который данный пользователь присвоит этому элементу, оценивается по рейтингам, данным этим пользователем другим близким элементам. И снова главное предположение состоит в том, что пользователю, положительно оценившему несколько элементов в прошлом, вероятно, понравятся элементы, оцененные аналогично этим прошлым выборам многими другими пользователями.

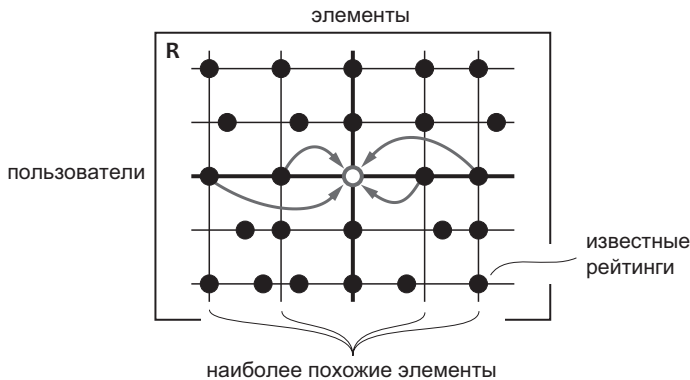


Рис. 5.9. Совместная фильтрация по близости элементов

Оба подхода к совместной фильтрации, по близости пользователей и элементов, требуют определения мер сходства между пользователями или элементами и некоторого метода усреднения рейтинга. Несмотря на структурное сходство подходов, они могут использовать разные меры, и для каждого подхода существует множество разных вариантов формул сходства и усреднения рейтинга. Мы подробно рассмотрим эти детали в следующих разделах.

5.7.1. Совместная фильтрация по близости пользователей

Два основных этапа совместной фильтрации на основе близости — это выбор пользователей или элементов, которых можно считать близкими, и прогнозирование рейтинга путем усреднения рейтингов соседей. Для подхода с пользователями это означает, что необходимо определить две ключевые функции: меру сходства пользователей и формулу усреднения рейтинга. В научной литературе и промышленных отчетах описано много разных вариантов этих функций, поэтому сначала опишем один из самых простых и известных вариантов, а затем обсудим возможные разновидности и улучшения.

Чтобы определить меру сходства, рассмотрим двух пользователей, u и v , оценивших элементы I_u и I_v соответственно. Множество элементов, оцененных обоими пользователями, определяется как пересечение:

$$I_{uv} = I_u \cap I_v. \quad (5.38)$$

На основе этого множества элементов можно определить меру сходства. Чаще всего на практике используется коэффициент корреляции Пирсона, который вычисляется следующим образом [Herlocker et al., 1999]:

$$\text{sim}(u, v) = \frac{\sum_{i \in I_{uv}} (r_{ui} - \mu_u)(r_{vi} - \mu_v)}{\sqrt{\sum_{i \in I_{uv}} (r_{ui} - \mu_u)^2} \sqrt{\sum_{i \in I_{uv}} (r_{vi} - \mu_v)^2}}, \quad (5.39)$$

где μ_u и μ_v — средние рейтинги пользователей:

$$\mu_u = \frac{1}{|I_{uv}|} \sum_{i \in I_{uv}} r_{ui}. \quad (5.40)$$

Обратите внимание, что формула 5.40 вычисляет средний рейтинг пользователей по набору общих элементов I_{uv} , как того требует определение коэффициента корреляции Пирсона. Следовательно, это значение не является постоянным для данного пользователя u , но уникально для каждой пары пользователей. На практике, однако, довольно часто используется глобальный средний рейтинг для пользователя u , вычисленный по всем элементам I_u , оцененным этим пользователем [Aggarwal, 2016].

Мера сходства позволяет выявить k пользователей, наиболее похожих на целевого пользователя u . Размер выборки k является параметром алгоритма рекомендаций. Поскольку наша цель — спрогнозировать рейтинг пользователя u и элемента i усреднением рейтингов, присвоенных этому элементу другими пользователями, мы выбираем не просто k самых похожих пользователей, но k самых похожих пользователей, оценивших элемент i . Обозначим это множество соседей как S_{ui}^k . Оно может включать меньше k пользователей, если матрица рейтингов не содержит достаточного количества оценок для элемента i или достаточного количества пользователей, похожих на u , с общими оцененными элементами. Далее рейтинг можно оценить как средневзвешенное рейтингов похожих пользователей:

$$\hat{r}_{ui} = \mu_u + \frac{\sum_{v \in S_{ui}^k} \text{sim}(u, v) \cdot (r_{vi} - \mu_v)}{\sum_{v \in S_{ui}^k} |\text{sim}(u, v)|}. \quad (5.41)$$

Формула 5.41 использует идею отделения смещений пользователей от сигнала взаимодействия, которую мы рассмотрели выше, в разделе, посвященном базовым оценкам. На первом этапе глобальные средние смещений пользователей, μ_u и μ_v , вычитаются из исходных рейтингов, затем вычисляется сигнал взаимодействия как произведение мер сходства и рейтингов, центрированных по среднему, и, наконец, смещение пользователя μ_u добавляется обратно для учета предпочтений целевого пользователя.

Рассмотрим несколько вариантов формул 5.39 и 5.41. Большинство из них являются эвристическими коррекциями, которые могут помочь улучшить точность

оценок и вычислительную стабильность в практических приложениях [Su and Khoshgoftaar, 2009; Breese et al., 1998]. Сначала рассмотрим несколько вариантов измерения сходства.

БАЗОВЫЕ ФУНКЦИИ СХОДСТВА. Коэффициент корреляции Пирсона, как известно, является хорошим способом измерения сходства, но также можно использовать другие метрики, включая косинусное сходство, коэффициент ранговой корреляции Спирмена и среднеквадратическое отклонение. Косинусное сходство между двумя пользователями, например, можно определить следующим образом:

$$\text{sim}(u, v) = \frac{\sum_{i \in I_{uv}} r_{ui} \cdot r_{vi}}{\sqrt{\sum_{i \in I_{uv}} r_{ui}^2} \sqrt{\sum_{i \in I_{uv}} r_{vi}^2}}. \quad (5.42)$$

Эти альтернативы обычно считаются слабее корреляции Пирсона, и, как мы увидим ниже, иногда эффективнее вычислять коэффициенты сходства с помощью регрессионного анализа, а не эвристического выбора и настройки функций подобия.

ВЗВЕШЕННОЕ СХОДСТВО. Мера сходства вычисляется только на основе элементов, оцененных обоими пользователями. Надежность этой оценки зависит от количества общих элементов, оцененных пользователями, поэтому часто полезно корректировать меру сходства в соответствии с поддержкой элементов (число общих рейтингов, которые имеют два элемента) [Koren, 2008]:

$$\text{sim}'(u, v) = \frac{|I_{uv}|}{|I_{uv}| + \lambda} \cdot \text{sim}(u, v). \quad (5.43)$$

Увеличивая параметр регуляризации, можно уменьшить ненадежные коэффициенты сходства с низкой поддержкой.

ОБРАТНАЯ ЧАСТОТА ПОЛЬЗОВАТЕЛЯ. Коэффициент корреляции Пирсона, как и многие другие стандартные метрики сходства, одинаково интерпретируют все элементы в множестве I_{uv} . Такой подход определенно не самый оптимальный, потому что некоторые элементы могут быть более показательными, чем другие. Например, если два пользователя положительно оценивают какой-то редкий или нишевый элемент, это, вероятно, в большей мере свидетельствует об их близости, чем если бы им обоим нравился очень популярный элемент. Основой для этой идеи послужила метрика подобия текста TF×IDF, в которой каждое слово взвешивается пропорционально его обратной частоте. Более формально обратную частоту пользователя (Inverse User Frequency, IUF) для элемента i можно определить как

$$w_i = \log \left(\frac{m}{|U_i|} \right), \quad (5.44)$$

где m — общее число пользователей, а $|U_i|$ — число пользователей, оценивших элемент i . Этот вес можно вставить в формулу коэффициента корреляции Пирсона:

$$\frac{\sum_{i \in I_{uv}} w_i (r_{ui} - \mu_u)(r_{vi} - \mu_v)}{\sqrt{\sum_{i \in I_{uv}} w_i (r_{ui} - \mu_u)^2} \sqrt{\sum_{i \in I_{uv}} w_i (r_{vi} - \mu_v)^2}}. \quad (5.45)$$

РЕЙТИНГИ ПО УМОЛЧАНИЮ. Метрика сходства для пользователей u и v обычно рассчитывается на пересечении оцененных элементов $I_u \cap I_v$. Учитывая разреженность матрицы рейтингов, это пересечение обычно мало, что снижает надежность оценок. Альтернативой является вычисление сходства на объединении оцененных элементов $I_u \cup I_v$, а не на пересечении, с добавлением некоторого нейтрального рейтинга по умолчанию для элементов, оцененных только одним из пользователей.

Функция прогнозирования рейтинга 5.41 смешивает известные центрированные значения рейтингов, используя в качестве весов оценки подобию. Есть несколько альтернативных вариантов этой функции, отличающихся логикой центрирования и взвешивания:

СТАНДАРТНАЯ ОЦЕНКА (Z-ОЦЕНКА). В статистике стандартной оценкой точки данных x (также известной как z-оценка) является отклонение этой точки от среднего значения, измеренное в стандартных отклонениях:

$$z(x) = \frac{x - \mu}{\sigma}. \quad (5.46)$$

Стандартную оценку можно рассматривать как нормализацию центрированных значений по стандартному отклонению. Мы можем использовать стандартные оценки как альтернативу центрированию в формуле прогнозирования рейтинга. Сначала вычисляется стандартная оценка рейтинга в контексте данного пользователя как

$$z(r_{ui}) = \frac{r_{ui} - \mu_u}{\sigma_u}, \quad (5.47)$$

где σ_u — стандартное отклонение известных пользовательских рейтингов:

$$\sigma_u = \sqrt{\frac{\sum_{i \in I_u} (r_{ui} - \mu_u)^2}{|I_u| - 1}}. \quad (5.48)$$

Теперь функцию прогнозирования рейтингов можно переопределить с использованием стандартных оценок вместо центрированных рейтингов:

$$\hat{r}_{ui} = \mu_u + \sigma_u \cdot \frac{\sum_{v \in S_{ui}^k} \text{sim}(u, v) \cdot z(r_{vi})}{\sum_{v \in S_{ui}^k} |\text{sim}(u, v)|}. \quad (5.49)$$

По аналогии с центрированием, исходные рейтинги сначала преобразуются с применением формулы 5.47, а затем выполняется обратное преобразование умножением результата на стандартное отклонение пользовательских рейтингов σ_u и обратным прибавлением среднего μ_u . Подход использования стандартной оценки усиливает рейтинги, полученные от пользователей с низкой дисперсией оценок, и уменьшает вес рейтингов пользователей с высокой дисперсией.

ЦЕНТРИРОВАНИЕ ПО БАЗОВОЙ ВЕЛИЧИНЕ. Другой альтернативой центрированию по среднему является центрирование по базовой величине, использующее базовые оценки, вычисляемые по формуле 5.32 [Koren, 2008]. Формулу прогнозирования с центрированием по базовой величине можно определить следующим образом:

$$\hat{r}_{ui} = b_{ui} + \frac{\sum_{v \in S_{ui}^k} \text{sim}(u, v) \cdot (r_{vi} - b_{ui})}{\sum_{v \in S_{ui}^k} |\text{sim}(u, v)|}. \quad (5.50)$$

УСИЛЕНИЕ. Стандартная функция прогнозирования рейтинга использует в качестве весов оценки сходства. К этим исходным оценкам можно применить некоторые нелинейные преобразования для усиления определенных рейтингов. Например, следующее преобразование усиливает высокие оценки сходства для параметра усиления $p > 1$:

$$\text{sim}'(u, v) = \text{sim}(u, v) \cdot |\text{sim}^{p-1}(u, v)|. \quad (5.51)$$

ВЫБОР СОСЕДЕЙ. Качество рекомендаций обычно зависит от количества k пользователей, включенных в число соседей. Некоторые исследования показывают, что точность предсказания рейтингов может монотонно увеличиваться с увеличением числа соседей при условии использования продвинутой, правильно спроектированной и тщательно настроенной модели прогнозирования [Koren, 2008]. Конечно, постепенное улучшение за счет увеличения числа соседей постепенно сходит на нет и после определенного момента становится незначительным. Однако некоторые другие исследования показывают, что слишком большое число соседей может негативно влиять на точность рекомендаций, сделанных базовыми методами на основе близости, из-за шума, который приносят соседи с низким сходством [Herlocker et al., 1999; Bellogín et al., 2014]. Количество соседей k можно ограничить эмпирически, некоторым постоянным значением, или порогом подобия, который отфильтровывает соседей с небольшими оценками сходства.

ПРИМЕР 5.3

Кратко проиллюстрируем подход к совместной фильтрации по близости пользователей на численном примере, используя данные из табл. 5.3. Прежде всего оценим попарное сходство пользователей по формуле 5.39. В результате получится следующая матрица сходства:

$$\begin{array}{l}
 \text{User 1} \quad \text{User 2} \quad \text{User 3} \quad \text{User 4} \quad \text{User 5} \quad \text{User 6} \\
 \begin{array}{l}
 \text{Пользователь 1} \\
 \text{Пользователь 2} \\
 \text{Пользователь 3} \\
 \text{Пользователь 4} \\
 \text{Пользователь 5} \\
 \text{Пользователь 6}
 \end{array}
 \begin{bmatrix}
 1,00 & 0,87 & 0,94 & -0,79 & -0,59 & -0,78 \\
 0,87 & 1,00 & 0,87 & -0,84 & -0,81 & -0,88 \\
 0,94 & 0,87 & 1,00 & -0,93 & -0,87 & -0,91 \\
 -0,79 & -0,84 & -0,93 & 1,00 & 0,85 & 0,95 \\
 -0,59 & -0,81 & -0,87 & 0,85 & 1,00 & 0,94 \\
 -0,78 & -0,88 & -0,91 & 0,95 & 0,94 & 1,00
 \end{bmatrix}
 \end{array} \quad (5.52)$$

Как видите, первые три пользователя положительно коррелируют друг с другом и отрицательно коррелируют с последними тремя. Матрица сходства позволяет отыскать k наиболее похожих пользователей для заданного целевого пользователя и смешать их рейтинги, чтобы получить прогноз. Например, спрогнозируем недостающий рейтинг для пользователя 1 и фильма *The Godfather*, задав количество соседей k равным 2. Самыми похожими соседями в данном случае являются пользователи 3 и 2, присвоившие фильму *The Godfather* рейтинги 5 и 3 соответственно. Применяв формулу прогнозирования рейтинга 5.41, получим следующую оценку:

$$\begin{aligned}
 \hat{r}_{13} &= \mu_1 + \frac{\text{sim}(1,3) \cdot (r_{33} - \mu_3) + \text{sim}(1,2) \cdot (r_{23} - \mu_2)}{|\text{sim}(1,3)| + |\text{sim}(1,2)|}, \\
 &= 2,60 + \frac{0,94 \cdot (5 - 4,00) + 0,87 \cdot (3 - 2,20)}{0,94 + 0,87} = 3,50.
 \end{aligned} \quad (5.53)$$

Повторяя этот процесс для всех отсутствующих рейтингов, получим результаты, показанные в табл. 5.5. Обратите внимание, что эти оценки выглядят более понятными и точными, чем базовые оценки в табл. 5.4.

На практике методы рекомендаций по пользователям могут столкнуться с проблемами масштабируемости, по мере увеличения числа пользователей системы до десятков и сотен миллионов. Если количество ближайших соседей для целевого пользователя определяется в момент запроса рекомендаций, необходимо

вычислить метрики сходства между целевым пользователем и всеми другими пользователями системы. Если количество ближайших соседей определяется заранее, объем вычислений можно выразить как квадратичную функцию от числа пользователей. Кроме того, профиль целевого пользователя (например, история просмотров в текущем веб-сеансе) может быть недоступен заранее. Один из возможных способов обойти это ограничение — заменить сходство пользователей сходством элементов, как рассказывается в следующем разделе.

Таблица 5.5. Пример прогнозирования рейтингов с использованием алгоритма совместной фильтрации по пользователям

	Forrest Gump	Titanic	The Godfather	Batman	The Matrix	Alien
Пользователь 1	5	4	[3,50]	1	2	1
Пользователь 2	4	[3,40]	3	1	1	2
Пользователь 3	[6,11]	5	5	[2,59]	3	3
Пользователь 4	2	[2,64]	1	4	5	4
Пользователь 5	2	2	2	[3,62]	4	[3,61]
Пользователь 6	1	2	1	[3,76]	5	4

5.7.2. Совместная фильтрация по близости элементов

Основная идея подхода к совместной фильтрации по близости элементов заключается в выборе для рекомендаций элементов, похожих на элементы, прежде положительно оцененные целевым пользователем, путем вычисления меры сходства между элементами на основе известных рейтингов, присвоенных другими пользователями. Он напоминает подход рекомендаций по содержанию, в том смысле что рекомендации выбираются на основе сходства элементов, хотя метрики сходства имеют совершенно разную природу. В то же время этот подход структурно похож на совместную фильтрацию по близости пользователей, поскольку оба метода основаны на понятии ближайшего окружения и, следовательно, имеют одну и ту же алгоритмическую основу [Linden et al., 2003; Sarwar et al., 2001].

Чтобы предсказать рейтинг, который пользователь u присвоит элементу i , система рекомендаций по близости элементов сначала определяет ближайших соседей элемента i , то есть множество k наиболее похожих элементов. Чтобы вычислить

меру сходства между двумя элементами, i и j , обозначим множество пользователей, оценивших элемент i , как U_i , множество пользователей, оценивших элемент j , как U_j , и множество пользователей, оценивших оба элемента, как

$$U_{ij} = U_i \cap U_j. \quad (5.54)$$

Теперь сходство можно измерить как коэффициент корреляции Пирсона между векторами с общими рейтингами:

$$\text{sim}(i, j) = \frac{\sum_{u \in U_{ij}} (r_{ui} - \mu_i)(r_{uj} - \mu_j)}{\sqrt{\sum_{u \in U_{ij}} (r_{ui} - \mu_i)^2} \sqrt{\sum_{u \in U_{ij}} (r_{uj} - \mu_j)^2}}, \quad (5.55)$$

где μ_i и μ_j — средние рейтинги для элементов i и j соответственно. Эта формула совпадает с корреляцией Пирсона для пользователей в уравнении 5.39; единственное отличие — пользователи (строки) заменяются элементами (столбцами). Все элементы, оцененные пользователем u , можно теперь отсортировать по сходству с данным элементом i и из этого списка выбрать k наиболее похожих элементов. Обозначим множество соседей, ближайших к элементу i , как Q_{ui}^k . Обратите внимание, что множество включает только элементы, оцененные целевым пользователем u , а не все наиболее похожие элементы в каталоге, поэтому с увеличением k множество Q_{ui}^k сходится к I_u . После этого рейтинги можно предсказать как средневзвешенные рейтинги по k наиболее похожим элементам на основе центрированных по среднему рейтингов:

$$\hat{r}_{ui} = \mu_i + \frac{\sum_{j \in Q_{ui}^k} \text{sim}(i, j) \cdot (r_{uj} - \mu_j)}{\sum_{j \in Q_{ui}^k} |\text{sim}(i, j)|}. \quad (5.56)$$

По аналогии с подходом, основанным на сходстве пользователей, формулы 5.55 и 5.56 определяют лишь базовые оценки, которые можно скорректировать и улучшить с помощью разных приемов, которые мы рассматривали выше, таких как взвешенное сходство и усиление. Большинство из этих приемов применимо к обоим методам, по сходству пользователей и элементов. Например, для повышения точности прогнозов можно использовать входные рейтинги, центрированные по базовому уровню, а не по среднему [Koren, 2008]:

$$\hat{r}_{ui} = b_{ui} + \frac{\sum_{j \in Q_{ui}^k} \text{sim}(i, j) \cdot (r_{uj} - b_{uj})}{\sum_{j \in Q_{ui}^k} |\text{sim}(i, j)|}. \quad (5.57)$$

5.7.3. Сравнение методов на основе близости пользователей и элементов

Подход на основе близости элементов был предложен спустя годы после появления первых методов на основе близости пользователей, но он быстро приобрел популярность из-за лучшей масштабируемости и вычислительной эффективности [Linden et al., 2003; Koren and Bell, 2011]. Одно из ключевых преимуществ заключается в том, что общее количество элементов m в системе часто достаточно мало, чтобы можно было предварительно рассчитать и сохранить матрицу сходств элементов $m \times m$ и быстрее найти лучшие k рекомендаций для данного профиля пользователя. Это гарантирует лучшую масштабируемость архитектуры рекомендательной системы: тяжелые вычисления, необходимые для создания матрицы сходств, выполняются в фоновом режиме, и эта матрица используется для выбора рекомендаций в режиме реального времени. Конечно, эту стратегию можно также применить к методам на основе близости пользователей, но в рекомендательных системах с большим числом пользователей она может получиться очень дорогостоящей или совершенно непрактичной. Наконец, некоторые исследования показали, что методы на основе близости элементов систематически превосходят методы на основе близости пользователей с точки зрения точности прогнозирования для некоторых важных наборов данных, таких как данные Netflix [Bell and Koren, 2007].

В то же время важно отметить, что подходы на основе близости пользователей способны распознавать определенные взаимосвязи, не определяемые методами на основе близости элементов [Koren and Bell, 2011]. Напомню, что подход на основе близости элементов предсказывает рейтинг r_{ui} , опираясь на рейтинги, которые пользователь u присвоил элементам, похожим на элемент i . Этот прогноз едва ли будет точным, если ни один из элементов, оцененных пользователем, не похож на i . С другой стороны, подход на основе близости пользователей может выявить пользователей, похожих на u и оценивших i , и дать более надежный прогноз рейтинга. Как мы увидим далее, некоторые улучшенные методы рекомендаций сочетают модели на основе близости элементов и пользователей, чтобы воспользоваться преимуществами обоих методов.

Отношение между количеством пользователей и элементов является одним из ключевых факторов, влияющих на выбор того или иного подхода. В розничной торговле, например, часто предпочтительнее выглядит подход на основе близости элементов, потому что число элементов меньше числа пользователей. Однако в некоторых областях число элементов может превышать число пользователей. Например, в системе рекомендаций статей для исследователей может оказаться выгоднее использовать решение на основе близости пользователей, потому что общее количество всех опубликованных научных статей достигает многих сотен миллионов, тогда как численность исследовательского сообщества, пользующегося системой, намного менее многочисленно [Jack et al., 2016].

5.7.4. Методы на основе близости как задача регрессии

Методы на основе близости, рассматривавшиеся в предыдущих разделах, полагаются на эвристическую функцию прогнозирования рейтинга, которая оценивает неизвестные рейтинги как взвешенные средние известных рейтингов. Чтобы сделать утверждение о взвешенных средних более явным, отметим, что функции 5.41 и 5.56 прогнозирования рейтингов на основе близости пользователей и элементов по существу имеют следующие формы:

$$\hat{r}_{ui} = \sum_{v \in S_{ui}^k} w_{uv} \cdot r_{vi} \quad (\text{на основе близости пользователя}), \quad (5.58)$$

$$\hat{r}_{ui} = \sum_{j \in Q_{ui}^k} w_{ij} \cdot r_{uj} \quad (\text{на основе близости элемента}), \quad (5.59)$$

где w_{uv} и w_{ij} — весовые коэффициенты, пропорциональные сходствам пользователей и элементов соответственно. Другими словами, весовые коэффициенты w являются коэффициентами интерполяции. Это соображение вполне естественно приводит к вопросу определения оптимальных весов с использованием регрессионного анализа вместо эвристических весов на основе сходства. Регрессионный анализ можно применить к обеим моделям, как на основе пользователей, так и на основе элементов, а также к гибридным методам, объединяющим эти две модели, поэтому начнем с самого, пожалуй, практичного подхода, основанного на элементах, а затем обсудим альтернативные варианты [Bell and Koren, 2007].

5.7.4.1. Регрессия по элементам

Методы на основе близости элементов предсказывают рейтинги для элемента i , усреднения рейтингов похожих элементов, согласно выражению 5.59. Входные рейтинги r_{uj} можно взять непосредственно из исходной матрицы рейтингов, также матрицу можно предварительно обработать, центрировав рейтинги вычитанием глобального среднего, среднего для элемента или базового прогноза. В случае с центрированными входными рейтингами выходной рейтинг \hat{r}_{ui} также оказывается центрированным, поэтому в конце необходимо обратно прибавить глобальное среднее, среднее для элемента или базовый прогноз.

Чтобы решить регрессионную задачу для коэффициентов интерполяции рейтинга, рассмотрим гипотетический случай, когда матрица рейтингов настолько плотная, что все пользователи, кроме u оценили элемент i и всех его соседей Q_{ui}^k , как показано на рис. 5.10.

В этом случае оптимальные коэффициенты интерполяции для элемента i можно определить, решив следующую задачу наименьших квадратов (для каждого элемента отдельно):

$$\min_w \sum_{v \neq u} (r_{vi} - \hat{r}_{vi})^2. \quad (5.60)$$

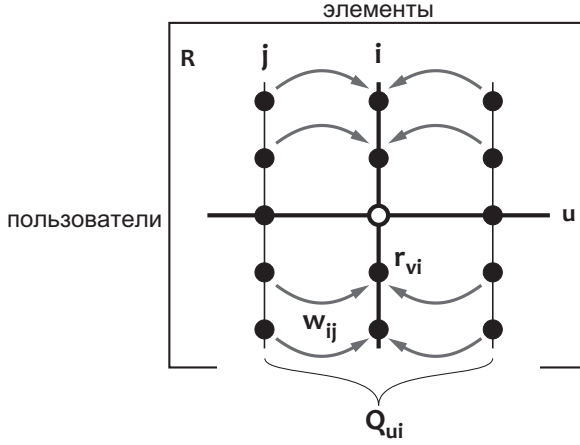


Рис. 5.10. Регрессия на основе близости элементов

Вставив в эту формулу функцию предсказания рейтингов 5.59, получим

$$\min_w \sum_{v \neq u} \left(r_{vi} - \sum_{j \in Q_{vi}^k} w_{ij} \cdot r_{vj} \right)^2. \quad (5.61)$$

Переупорядочив члены, можно переписать эту задачу в векторной форме:

$$\min_w \mathbf{r}^T \mathbf{r} - 2\mathbf{b}^T \mathbf{w} + \mathbf{w}^T \mathbf{A} \mathbf{w}, \quad (5.62)$$

где \mathbf{A} — матрица $k \times k$, которая определяется как

$$A_{jh} = \sum_{v \neq u} r_{vj} r_{vh}, \quad (5.63)$$

\mathbf{b} — k -мерный вектор, который определяется как

$$b_j = \sum_{v \neq u} r_{vj} r_{vi}, \quad (5.64)$$

и $\mathbf{r}^T \mathbf{r}$ — постоянный член относительно \mathbf{w} :

$$\mathbf{r}^T \mathbf{r} = \sum_{v \neq u} r_{vi}^2. \quad (5.65)$$

Если взять градиент квадратичной формы 5.62 относительно \mathbf{w} и приравнять его к нулю, получится следующая система линейных уравнений:

$$\mathbf{A}\mathbf{w} = \mathbf{b}. \quad (5.66)$$

В более реалистичном случае, с разреженной матрицей рейтингов, можно ожидать, что лишь несколько пользователей оценят элемент i и всех его ближайших соседей Q_{ui}^k . Следовательно, оценки соответствующих элементов \mathbf{A} и \mathbf{b} могут оказаться более или менее надежными, в зависимости от количества известных рейтингов. Мы можем объяснить это, используя метод взвешенного сходства, который рассматривался выше, и сократив оценки соответствующей поддержкой:

$$\begin{aligned} A_{jh} &= \frac{1}{|U_{jh}|} \sum_{v \in U_{jh}} r_{vj} r_{vh}, \\ b_j &= \frac{1}{|U_{ij}|} \sum_{v \in U_{ij}} r_{vj} r_{vi}, \end{aligned} \quad (5.67)$$

где U_{ij} — множество пользователей, оценивших оба элемента, i и j . Стоит отметить, что можно заранее вычислить и хранить все возможные элементы матрицы \mathbf{A} , то есть вычислить $m \times m$ элементов корреляционной матрицы в соответствии с выражением 5.63 для всех значений $1 \leq j, k \leq m$, а затем использовать эти значения, чтобы быстро собрать $k \times k$ матрицу \mathbf{A} и вектор \mathbf{b} для данного элемента и целевого пользователя.

Одним из способов вычисления оптимальных весов является численное решение уравнения 5.66 путем инвертирования матрицы \mathbf{A} , но это не единственный возможный вариант. Альтернативный подход заключается в непосредственном решении задачи 5.62 методом градиентного спуска или другого универсального метода оптимизации. Преимущество данного подхода — в возможности добавления дополнительных ограничений и переменных. Например, было подмечено: точность прогнозирования немного улучшается, если весовые коэффициенты \mathbf{w} ограничить неотрицательными значениями [Bell and Koren, 2007]:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \mathbf{w}^T \mathbf{A} \mathbf{w} - 2 \mathbf{b}^T \mathbf{w} \\ \text{при условии} \quad & \mathbf{w} \geq 0. \end{aligned} \quad (5.68)$$

Еще лучшие результаты можно получить, добавив больше переменных в базовую формулу 5.59 прогнозирования рейтинга и совместно оптимизировав их. Например, было показано, что следующее расширение формулы прогнозирования рейтинга является хорошим практическим выбором для базовых оценок [Koren and Bell, 2011]:

$$\hat{r}_{ui} = \mu + b_u + b_i + \sum_{j \in Q_{ui}^k} (w_{ij}(r_{uj} - b_{ui}) + c_{ij}), \quad (5.69)$$

где μ — глобальный средний рейтинг, b_{ij} — базовые оценки и b_u, b_i, w_{ij} и c_{ij} — переменные для оптимизации. Это выражение можно вставить в задачу наименьших квадратов 5.60 и оптимизировать по отношению к b_u, b_i, w_{ij} и c_{ij} методом градиентного спуска. В этом случае можно не ограничивать количество ближайших соседей k элементами и использовать весь набор I_u вместо Q_{ui}^k .

ПРИМЕР 5.4

Завершим обзор регрессии по элементам численным примером, в котором используется матрица рейтингов фильмов из табл. 5.3. Будем работать с рейтингами, центрированными по среднему, поэтому сначала обработаем исходную матрицу, вычтя среднее по столбцу (то есть средний рейтинг элемента) из каждого элемента и получим результат, представленный в табл. 5.6.

Таблица 5.6. Пример рейтингов, центрированных по среднему

	Forrest Gump	Titanic	The Godfather	Batman	The Matrix	Alien
Среднее	2,80	3,25	2,40	2,00	3,33	2,80
Пользователь 1	2,20	0,75	–	–1,00	–1,33	–1,80
Пользователь 2	1,20	–	0,60	–1,00	–2,33	–0,80
Пользователь 3	–	1,75	2,60	–	–0,33	0,20
Пользователь 4	–0,80	–	–1,40	2,00	1,66	1,20
Пользователь 5	–0,80	–1,25	–0,40	–	0,66	–
Пользователь 6	–1,80	–1,25	–1,40	–	1,66	1,20

Далее вычислим матрицу корреляций A элементов по формуле 5.63 для всех значений $1 \leq j, k \leq m$ рейтингов из табл. 5.6:

$$A = \begin{bmatrix} 10,80 & 4,90 & 4,68 & -5,00 & -10,60 & -8,04 \\ 4,90 & 6,75 & 6,80 & -0,75 & -4,50 & -2,50 \\ 4,68 & 6,80 & 11,20 & -3,40 & -7,20 & -3,32 \\ -5,00 & -0,75 & -3,40 & 6,00 & 7,00 & 5,00 \\ -10,6 & -4,50 & -7,20 & 7,00 & 13,33 & 8,20 \\ -8,04 & -2,50 & -3,32 & 5,00 & 8,20 & 6,80 \end{bmatrix}. \quad (5.70)$$

Эту предварительно вычисленную матрицу можно использовать, чтобы быстро собрать систему уравнений, описанную выражением 5.66, для заданного целевого пользователя и элемента. Например, с количеством ближайших соседей $k=2$ предскажем рейтинг, который пользователь 1 присвоит фильму *The Godfather*, усреднив рейтинги фильмов *Titanic* и *Forrest Gump*, которые являются двумя ближайшими соседями для *The Godfather* с точки зрения сходства Пирсона. Соответственно используем значения корреляций для этих трех фильмов (выделены в матрице 5.70), чтобы построить следующее уравнение для коэффициентов интерполяции:

$$\begin{bmatrix} 6,75 & 4,90 \\ 4,90 & 10,80 \end{bmatrix} \begin{bmatrix} w_{32} \\ w_{31} \end{bmatrix} = \begin{bmatrix} 6,80 \\ 4,68 \end{bmatrix}. \quad (5.71)$$

Решив это уравнение, получим весовые коэффициенты $w_{32} = 1,033$ для *Titanic* и $w_{31} = 0,035$ для *Forrest Gump*. Теперь можно предсказать рейтинг как

$$\begin{aligned} \hat{r}_{13} &= \mu_3 + w_{32}r_{12} + w_{31}r_{11} = \\ &= 2,40 + 1,033 \times 0,75 - 0,035 \times 2,20 = 3,09, \end{aligned} \quad (5.72)$$

где μ_3 — средний рейтинг фильма *The Godfather*, а входные рейтинги r_{12} и r_{11} взяты из табл. 5.6. Повторяя этот процесс для всех неизвестных оценок, получим окончательные результаты, представленные в табл. 5.7.

Таблица 5.7. Пример рейтингов, предсказанных методом регрессии по элементам

	Forrest Gump	Titanic	The Godfather	Batman	The Matrix	Alien
Пользователь 1	5	4	[3,09]	1	2	1
Пользователь 2	4	[3,83]	3	1	1	2
Пользователь 3	[4,02]	5	5	[1,98]	3	3
Пользователь 4	2	[2,34]	1	4	5	4
Пользователь 5	2	2	2	[2,28]	4	[3,15]
Пользователь 6	1	2	1	[2,94]	5	4

5.7.4.2. Регрессия по пользователям

Аппарат регрессионного анализа для методов регрессии по элементам, который мы только что разработали, можно также применить к моделям на основе пользователей. Входными данными для обработки в данном случае является матрица рейтингов пользователей, центрированная по средним значениям (среднее значение для строки вычитается из каждого элемента в этой строке), или базовые прогнозы. Задачу наименьших квадратов 5.60 для случая регрессии по пользователям можно определить как

$$\min_W \sum_{j \neq i} (r_{ij} - \hat{r}_{ij})^2. \quad (5.73)$$

Эту задачу нужно решить для каждого целевого пользователя u . Вставив формулу прогноза рейтинга 5.58 в предыдущее уравнение, получим:

$$\min_W \sum_{j \neq i} \left(r_{ij} - \sum_{v \in S_{ij}^k} w_{uv} \cdot r_{vj} \right)^2. \quad (5.74)$$

Оптимальные весовые коэффициенты w_{uv} можно определить с помощью тех же методов, что и в подходе на основе элементов — решить систему линейных уравнений или использовать универсальные методы оптимизации весов на основе функции затрат 5.74. В отличие от подхода на основе элементов, регрессия по пользователям наследует все преимущества и недостатки методов на основе пользователей, которые рассматривались выше. В частности, эти методы сложнее с вычислительной точки зрения, если пользователей окажется намного больше, чем элементов, потому что требуется предварительно вычислить матрицу $n \times n$ корреляций пользователей, а не матрицу $m \times m$ корреляций элементов.

5.7.4.3. Сплав моделей на основе близости пользователей и элементов

Одно из ключевых преимуществ регрессионного подхода — возможность расширения функции прогнозирования рейтингов новыми членами и переменными, которые можно оптимизировать вместе. Мы уже видели пример такого расширения в выражении 5.69, где добавили новые переменные в базовую модель на основе элементов, чтобы учесть систематическую смещенность оценок элементов и предвзятость пользователей. Это решение можно расширить еще больше и объединить модели на основе элементов и пользователей в одну функцию прогнозирования рейтинга:

$$\hat{r}_{ui} = \mu + b_u + b_i + \sum_{j \in Q_{ui}^k} \left(w_{ij}^{(item)} (r_{uj} - b_{uj}) + c_{ij} \right) + \sum_{v \in S_{ui}^k} w_{uv}^{(user)} (r_{vi} - b_{vi}), \quad (5.75)$$

где $w_{ij}^{(item)}$ и $w_{uv}^{(user)}$ — два разных множества весовых коэффициентов, которые требуется оптимизировать. Эта модель, по сути, является суммой центрированной версии функции 5.58 на основе пользователей и функции 5.69 на основе элементов [Aggarwal, 2016; Koren and Bell, 2011]. Полученную функцию прогнозирования рейтинга можно вставить в задачу наименьших квадратов определения ошибки прогнозирования и оптимизации по всем смещенным переменным весовым коэффициентам. Поскольку веса определяются по данным, количество пользователей и элементов может не ограничиваться ближайшими k элементами, и вместо Q_{ui}^k и S_{ui}^k могут использоваться множества I_u и U_i соответственно. Однако если множества ограничить конечным значением k , вычислительную сложность можно уменьшить за счет точности модели.

Комбинированная модель позволяет выявить взаимосвязи элемент/элемент и пользователь/пользователь (подробности смотрите в разделе 5.7.3) и, как следствие, получить все преимущества обоих подходов. Как было подмечено, комбинированные модели могут превосходить индивидуальные модели прогнозирования по пользователям и элементам на промышленных наборах данных [Koren and Bell, 2011]. Важно отметить, что регрессионный аппарат можно использовать не только для объединения решений на основе пользователей и элементов, но и для интеграции методов, основанных на близости, с совершенно другими моделями, включая те, которые мы обсудим в следующем разделе.

5.8. Совместная фильтрация на основе моделей

С точки зрения машинного обучения подход к совместной фильтрации на основе близости является слишком узким взглядом на проблему, поскольку фокусируется на оценках k ближайших соседей и не использует другие методы машинного обучения. Следовательно, рекомендательная система на основе близости наследует некоторые фундаментальные ограничения метода k ближайших соседей. Во-первых, производительность методов на основе близости может снизиться на разреженных данных, где элементы или пользователи имеют мало общих рейтингов, из-за чего рекомендации могут делаться на основе соседей, на самом деле не похожих на целевого пользователя или элемент. Кроме того, алгоритм k ближайших соседей опирается на попарное сравнение и откладывает вычисление рекомендаций до

момента их запроса, что затрудняет разделение вычисления на этапы предварительных и оперативных вычислений.

Альтернативный подход заключается в построении модели прогнозирования рейтинга с использованием более продвинутых методов машинного обучения с учителем и без учителя. Совместная фильтрация — это, по сути, задача восстановления матрицы, поэтому для ее решения можно использовать множество стандартных методов классификации и регрессии. Этот подход, известный как совместная фильтрация на основе моделей, обычно имеет несколько преимуществ перед методами на основе близости.

ТОЧНОСТЬ. Некоторые методы машинного обучения, такие как наивный байесовский классификатор, основаны на прочном теоретическом фундаменте, который позволяет прогнозировать рейтинг более точно, чем эвристические меры сходства, используемые рекомендательными системами на основе близости.

СТАБИЛЬНОСТЬ. Методы уменьшения размерности позволяют преобразовать разреженную матрицу оценок в более сжатое представление, что повышает стабильность предсказания рейтингов по неполным данным.

МАСШТАБИРУЕМОСТЬ. Методы машинного обучения часто состоят из этапов обучения и использования модели, которые помогают отделить предварительные вычисления от интерактивных запросов рекомендаций, тем самым улучшая масштабируемость системы.

Некоторые методы на основе моделей могут обеспечить все эти улучшения, тогда как другие достигают только некоторых из них. В остальной части этого раздела мы рассмотрим несколько важных методов, которые могут превосходить системы на основе близости или объединяться с алгоритмами на основе близости для создания гибридных решений.

5.8.1. Адаптация регрессионных моделей для предсказания рейтингов

В общем случае модели классификации и регрессии можно адаптировать к задаче прогнозирования рейтингов, если рассматривать известные рейтинги как признаки, а отсутствующие — как переменные отклика. Сначала рассмотрим гипотетический случай, когда в матрице отсутствует только один рейтинг, а все остальные известны. По аналогии с подходом на основе близости мы имеем две симметричные альтернативы, в зависимости от интерпретации столбцов и строк матрицы рейтингов. Первый вариант — рассматривать столбцы матрицы рейтинга как признаки, а строки как образцы данных. Модель классификации обучается

отдельно для каждого элемента i , когда i -й столбец рассматривается как отклик, а другие столбцы — как признаки, то есть рейтинг для данного элемента прогнозируется на основе рейтингов других элементов, как показано на рис. 5.11. Этот подход структурно напоминает методы на основе близости элементов. Вторая альтернатива заключается в том, чтобы в качестве признаков использовать строки матрицы рейтингов, а столбцы — как образцы данных, то есть модель классификации создается для каждого пользователя, а рейтинги целевого пользователя прогнозируются на основе рейтингов других пользователей. Этот подход можно рассматривать как задачу прогнозирования на основе пользователей.

Однако на практике матрица рейтингов часто очень разрежена, поэтому нельзя полагаться на то, что все рейтинги в обучающих образцах известны. Это серьезная проблема, которая может существенно влиять на качество прогнозов в зависимости от способа обработки отсутствующих значений. Существует несколько возможных путей решения этой проблемы:

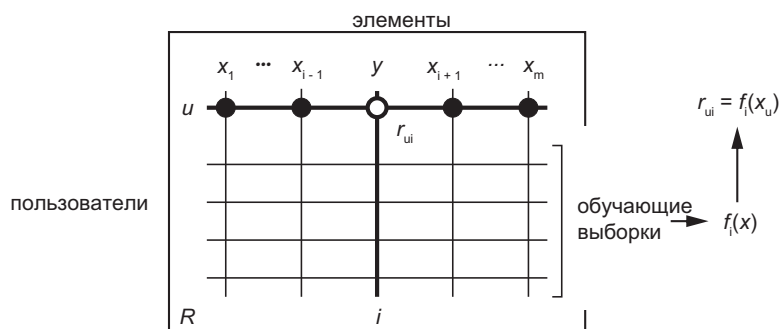


Рис. 5.11. Адаптация моделей регрессии или классификации для прогнозирования рейтингов на основе элементов. Известные рейтинги целевого пользователя u интерпретируются как признаки x_1, \dots, x_m , а прогнозируемый рейтинг интерпретируется как переменная отклика y . Модель регрессии или классификации $f_i(x)$ обучается для данного элемента i , используя других пользователей в качестве обучающих образцов, и применяется к вектору признаков x_u , который соответствует целевому пользователю

- Некоторые методы классификации можно непосредственно адаптировать для обработки отсутствующих значений. Примером такого решения является наивный байесовский классификатор, описанный в следующем разделе.
- В некоторых случаях отсутствующие рейтинги можно заменить нулями. Это в первую очередь относится к унарным матрицам рейтингов, где каждый элемент указывает, взаимодействовал ли пользователь с данным элементом или нет [Aggarwal, 2016]. Однако этот подход нельзя использовать повсеместно

для всех типов рейтингов, потому что вставка значений по умолчанию (ноль) приводит к смещению прогноза.

- Одним из наиболее универсальных решений является итеративный подход [Xia et al., 2006; Su et al., 2008]. Отсутствующие значения можно инициализировать некоторыми базовыми оценками, такими как средние значения строк или столбцов. Это дает полную матрицу рейтингов, которую можно использовать для обучения моделей классификации. Первоначально отсутствующие рейтинги можно оценить с помощью полученных классификаторов и подставить полученные значения в соответствующие элементы матрицы рейтингов. Это дает вторую полную рейтинговую матрицу, которую можно использовать для повторного обучения моделей, то есть этот процесс можно повторять итеративно до схождения.

Наконец, стоит отметить, что описанные выше методы можно смешивать друг с другом, а также с методами рекомендаций по содержанию и близости. Например, для инициализации отсутствующих рейтингов можно использовать наивный байесовский алгоритм совместной фильтрации, описанный в следующем разделе, а затем на основе этой полной матрицы рейтингов вычислить сходства Пирсона между пользователями или элементами, чтобы сгенерировать фактические рекомендации [Su et al., 2008].

5.8.2. Наивная байесовская совместная фильтрация

Наивный байесовский алгоритм совместной фильтрации пытается прогнозировать отсутствующие рейтинги путем оценки вероятностей всех возможных значений рейтинга (например, 1, 2, 3, 4 и 5 звезд) и выбора наиболее вероятного варианта [Miyahara and Pazzani, 2000; Su and Khoshgoftaar, 2006]. Как уже упоминалось в предыдущем разделе, наивный байесовский классификатор может быть ориентированным либо на пользователей, либо на элементы. Я предлагаю сосредоточиться на подходе, ориентированном на конкретные элементы, поскольку он, скорее всего, имеет большую практическую ценность по причинам, о которых говорилось в разделе 5.7.3. Решение, ориентированное на пользователей, можно построить почти таким же способом, поменяв местами пользователей и элементы, то есть строки и столбцы матрицы рейтингов.

Следуя элемент-ориентированному подходу, построим наивный байесовский классификатор для данного элемента i , чтобы предсказать рейтинг r_{ui} . Оценку рейтинга затем можно получить применением модели к множеству известных рейтингов пользователя u , которое мы обозначим как I_u . Если рейтинги являются категориальными переменными, принимающими значения из K возможных классов c_1, \dots, c_K ,

то задача прогнозирования состоит в том, чтобы найти наиболее вероятный класс рейтинга с учетом наблюдаемых рейтингов:

$$r_{ui} = \operatorname{argmax}_{c_k} \Pr(r_{ui} = c_k | I_u). \quad (5.76)$$

Чтобы получить вероятность определенного класса рейтинга с учетом наблюдаемых рейтингов, сначала применим правило Байеса для разложения этой вероятности:

$$\Pr(r_{ui} = c_k | I_u) = \frac{\Pr(c_k) \cdot \Pr(I_u | r_{ui} = c_k)}{\Pr(I_u)}, \quad (5.77)$$

где $\Pr(c_k)$ — априорная вероятность класса рейтинга c_k , а $\Pr(I_u | r_{ui} = c_k)$ — вероятность наблюдения известного рейтинга пользователя u , учитывая, что этот пользователь оценил элемент i как c_k . Вероятность наблюдаемых рейтингов $\Pr(I_u)$ в знаменателе можно игнорировать, так как она постоянна для всех классов и, следовательно, не влияет на выбор наиболее вероятного класса. Следующий шаг — применение наивного байесовского предположения для оценки вероятности наблюдаемых рейтингов. Согласно этому предположению, все наблюдаемые рейтинги считаются условно независимыми, поэтому вероятность можно разбить на произведение вероятностей индивидуальных рейтингов:

$$\Pr(I_u | r_{ui} = c_k) = \prod_{j \in I_u} \Pr(r_{uj} | r_{ui} = c_k). \quad (5.78)$$

Собрав промежуточные результаты вместе, получим итоговое выражение для прогнозирования рейтинга:

$$r_{ui} = \operatorname{argmax}_{c_k} \Pr(c_k) \cdot \prod_{j \in I_u} \Pr(r_{uj} | r_{ui} = c_k). \quad (5.79)$$

Последняя задача — оценить вероятности в уравнении 5.79 по данным. В контексте элемента i априорная вероятность класса рейтинга c_k оценивается как доля рейтингов для элемента i , которые равны c_k :

$$\Pr(c_k) = \frac{\sum_{v \in U_i} \mathbb{I}(r_{vi} = c_k)}{|U_i|}, \quad (5.80)$$

где U_i — это множество пользователей, оценивших элемент i , а $\mathbb{I}(x)$ — индикаторная функция, равная единице, если аргумент имеет истинное значение, и нулю в противном случае. Условную вероятность, что пользователь u присвоит элементу j

рейтинг r_{ij} с учетом того, что раньше он присвоил элементу i рейтинг c_k , можно оценить следующим образом:

$$\Pr(r_{ij} | r_{ui} = c_k) = \frac{\sum_{v \in U_i} \mathbb{I}(r_{vj} = r_{ij} \text{ AND } r_{vi} = c_k)}{\sum_{v \in U_i} \mathbb{I}(r_{vi} = c_k)}. \quad (5.81)$$

Числитель в выражении 5.81 равен числу пользователей, оценивших элемент j так же, как пользователь u , и в то же время оценивших элемент i как c_k . Знаменатель — это просто число пользователей, оценивших элемент i как c_k . Рассмотрим в качестве примера рис. 5.12. Оценим вероятность гипотезы, что r_{ui} — это 3 звезды. Допустим, что три пользователя присвоили элементу i 3 звезды и один из них присвоил элементу j такой же рейтинг, что и пользователь u (5 звезд), вероятность наблюдения известных рейтингов для элемента j , с учетом, что гипотеза верна, составит

$$\Pr(r_{uj} | r_{ui} = 3) = \frac{2}{3}. \quad (5.82)$$

На практике формула 5.81 часто корректируется с использованием алгоритма Лапласа, чтобы избежать нулевых количеств и сгладить оценки. Если принять, что $|C|$ — это общее число классов рейтинга, тогда сглаженная версия оценки правдоподобия будет выглядеть следующим образом:

$$\Pr(r_{ij} | r_{ui} = c_k) = \frac{\sum_{v \in U_i} \mathbb{I}(r_{vj} = r_{ij} \text{ AND } r_{vi} = c_k) + 1}{\sum_{v \in U_i} \mathbb{I}(r_{vi} = c_k) + |C|}. \quad (5.83)$$

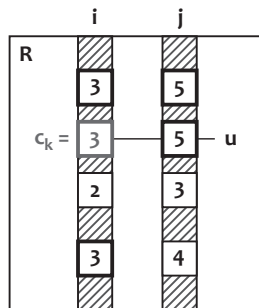


Рис. 5.12. Пример оценки правдоподобия в наивном байесовском алгоритме выбора рекомендаций

ПРИМЕР 5.5

Рассмотрим теперь полный числовой пример, используя стандартную матрицу рейтингов фильмов из табл. 5.3. У нас есть пять классов рейтингов, от 1 до 5 звезд, и входная матрица рейтингов выглядит следующим образом:

$$R = \begin{bmatrix} 5 & 4 & — & 1 & 2 & 1 \\ 4 & — & 3 & 1 & 1 & 2 \\ — & 5 & 5 & — & 3 & 3 \\ 2 & — & 1 & 4 & 5 & 4 \\ 2 & 2 & 2 & — & 4 & — \\ 1 & 2 & 1 & — & 5 & 4 \end{bmatrix}. \quad (5.84)$$

Возьмем в качестве примера отсутствующий рейтинг r_{13} и рассмотрим процедуру расчетов, которые необходимо выполнить, чтобы предсказать его величину с помощью наивного байесовского алгоритма, основанного на элементах. Первый шаг — оценка априорных вероятностей классов, согласно выражению 5.80. В результате получается следующий вектор вероятностей для классов от 1 до 5, согласно частотам в третьем столбце матрицы рейтингов:

$$\Pr(c_k) \begin{matrix} c_1 & c_2 & c_3 & c_4 & c_5 \\ [2/5 & 1/5 & 1/5 & 0 & 1/5]. \end{matrix} \quad (5.85)$$

Далее оценим условные вероятности, согласно выражению 5.83, относительно целевого элемента $i = 3$, для всех классов c_k и всех элементов $j \neq 3$. В результате получаем следующую матрицу вероятностей:

$$\begin{matrix} & c_1 & c_2 & c_3 & c_4 & c_5 \\ \begin{matrix} j=1 \\ j=2 \\ j=3 \\ j=4 \\ j=5 \\ j=6 \end{matrix} & \begin{bmatrix} 1/7 & 1/6 & 1/6 & 1/5 & 1/6 \\ 1/7 & 1/6 & 1/6 & 1/5 & 1/6 \\ — & — & — & — & — \\ 1/7 & 1/6 & 1/3 & 1/5 & 1/6 \\ 1/7 & 1/6 & 1/6 & 1/5 & 1/6 \\ 1/7 & 1/6 & 1/6 & 1/5 & 1/6 \end{bmatrix} \end{matrix}. \quad (5.86)$$

Матрицу выше можно рассчитать по запросу для данного целевого элемента, или предварительно вычислить все $m \times m \times |C|$ значений для всех возможных комбинаций целевых элементов i , соседних элементов j и классов c_k . Умножая значения в матрицах 5.85 и 5.86 по столбцам, согласно выражению 5.79, получаем следующие вероятности классов:

$$\Pr(r_{13} = c_k | I_1) \begin{matrix} c_1 & c_2 & c_3 & c_4 & c_5 \\ [2/84035 & 1/38880 & 1/19440 & 0 & 1/38880]. \end{matrix} \quad (5.87)$$

Это означает, что рейтинг 3 является лучшей оценкой для r_{13} . Если повторить процесс для всех отсутствующих рейтингов, то мы получим результаты, представленные в табл. 5.8.

Таблица 5.8. Пример рейтингов, предсказанных наивным байесовским алгоритмом совместной фильтрации по элементам

	Forrest Gump	Titanic	The Godfather	Batman	The Matrix	Alien
Пользователь 1	5	4	[3]	1	2	1
Пользователь 2	4	[4]	3	1	1	2
Пользователь 3	[2]	5	5	[1]	3	3
Пользователь 4	2	[2]	1	4	5	4
Пользователь 5	2	2	2	[4]	4	[4]
Пользователь 6	1	2	1	[4]	5	4

Попутно покажем также, как наивный байесовский алгоритм можно связать с совместной фильтрацией на основе близости. Их структурное сходство можно сделать более очевидным, заменив произведение в уравнении 5.78 суммой логарифмов и вставив его в формулу 5.77 вероятности класса:

$$\Pr(r_{ui} = c_k | I_u) = \frac{\Pr(c_k)}{\Pr(I_u)} \cdot \sum_{j \in I_u} s_k(i, j), \quad (5.88)$$

где

$$s_k(i, j) = \log \Pr(r_{ij} | r_{ui} = c_k). \quad (5.89)$$

Обратите внимание, что $s_k(i, j)$ оценивается путем попарного сравнения рейтингов элементов i и j , то есть это значение можно интерпретировать как своеобразную меру сходства между этими двумя элементами. Этот результат можно сравнить с формулой 5.56 ближайших соседей на основе элементов, согласно которой оценки предсказываются с помощью эвристической метрики сходства элементов. Этот более точный фундамент дает некоторое преимущество наивному байесовскому

алгоритму перед основными методами на основе близости, и он может существенно превосходить их на некоторых наборах данных [Miyahara and Pazzani, 2000].

5.8.3. Модели скрытых факторов

В алгоритмах совместной фильтрации, обсуждавшихся до сих пор, большая часть вычислений выполняется на основе отдельных элементов матрицы рейтингов. Методы на основе близости оценивают отсутствующие рейтинги непосредственно по известным значениям в матрице рейтингов. Методы на основе моделей добавляют слой абстракции поверх матрицы рейтингов, создавая предиктивную модель, которая фиксирует определенные закономерности взаимоотношений между пользователями и элементами, но обучение модели по-прежнему сильно зависит от свойств матрицы рейтингов. Как следствие, эти методы совместной фильтрации обычно сталкиваются со следующими проблемами:

- Матрица рейтингов может содержать миллионы пользователей, миллионы элементов и миллиарды известных рейтингов, что создает серьезные проблемы вычислительной сложности и масштабируемости.
- Матрица рейтингов, как правило, очень разрежена (на практике может отсутствовать около 99 % рейтингов). Это влияет на вычислительную стабильность алгоритмов рекомендаций и приводит к недостоверным оценкам, когда у пользователя или элемента нет действительно похожих соседей. Эта проблема часто усугубляется тем, что большинство базовых алгоритмов ориентированы либо на пользователей, либо на элементы, что ограничивает их способность фиксировать все типы сходств и взаимоотношений, доступных в матрице рейтингов.
- Данные в матрице рейтингов обычно сильно коррелируют из-за сходств пользователей и элементов. Это означает, что сигналы, доступные в матрице рейтингов, не только разрежены, но и избыточны, что способствует обострению проблемы масштабируемости.

Приведенные выше соображения указывают на то, что исходная матрица рейтингов может быть не самым оптимальным представлением сигналов, и следует рассмотреть другие альтернативные представления, более подходящие для целей совместной фильтрации. Чтобы изучить эту идею, вернемся к исходной точке и немного поразмышляем о характере служб рекомендаций. По сути, службу рекомендаций можно рассматривать как алгоритм, предсказывающий рейтинги на основе некоторой меры сходства между пользователем и элементом:

$$\hat{r}_{ui} \sim \text{affinity}(u, i). \quad (5.90)$$

Один из способов определить эту меру сходства — использовать подход скрытых факторов и отобразить пользователей и элементы в точки в некотором k -мерном пространстве, чтобы каждый пользователь и каждый элемент были представлены k -мерным вектором:

$$\begin{aligned} u &\rightarrow p_u = (p_{u1}, \dots, p_{uk}), \\ i &\rightarrow q_i = (q_{i1}, \dots, q_{ik}). \end{aligned} \quad (5.91)$$

Векторы должны строиться так, чтобы соответствующие размерности \mathbf{p} и \mathbf{q} были сопоставимы друг с другом. Иначе говоря, каждое измерение можно рассматривать как признак или понятие, то есть p_{uj} является мерой близости пользователя u и понятия j , а q_{ij} , соответственно, является мерой элемента i и понятия j . На практике эти размерности часто интерпретируются как жанры, стили и прочие атрибуты, применимые одновременно к пользователям и элементам. Сходство между пользователем и элементом и, соответственно, рейтинг можно определить как произведение соответствующих векторов:

$$\hat{r}_{ui} = \mathbf{p}_u \cdot \mathbf{q}_i^T = \sum_{s=1}^k p_{us} q_{is}. \quad (5.92)$$

Поскольку каждый рейтинг можно разложить на произведение двух векторов, принадлежащих пространству понятий, которое не наблюдается непосредственно в исходной матрице рейтингов, \mathbf{p} и \mathbf{q} называются *скрытыми факторами*. Успех этого абстрактного подхода, конечно, полностью зависит от того, как именно определяются и конструируются скрытые факторы. Чтобы ответить на этот вопрос, заметим, что выражение 5.92 можно переписать в матричной форме следующим образом¹:

$$\hat{\mathbf{R}} = \mathbf{P} \cdot \mathbf{Q}^T, \quad (5.93)$$

где \mathbf{P} — матрица $n \times k$, собранная из векторов \mathbf{p} , а \mathbf{Q} — матрица $m \times k$, собранная из векторов \mathbf{q} , как показано на рис. 5.13. Основной целью системы совместной фильтрации обычно является минимизация ошибки прогнозирования рейтинга, что позволяет прямо определить задачу оптимизации относительно матриц скрытых факторов:

$$\min_{\mathbf{P}, \mathbf{Q}} \|\mathbf{R} - \hat{\mathbf{R}}\|^2 = \|\mathbf{R} - \mathbf{P} \cdot \mathbf{Q}^T\|^2. \quad (5.94)$$

¹ В математической литературе такие факторы часто обозначают, как U и V . Мы обозначим их как \mathbf{P} и \mathbf{Q} , чтобы избежать путаницы с индексом u , обозначающим пользователя, широко используемым в книге.

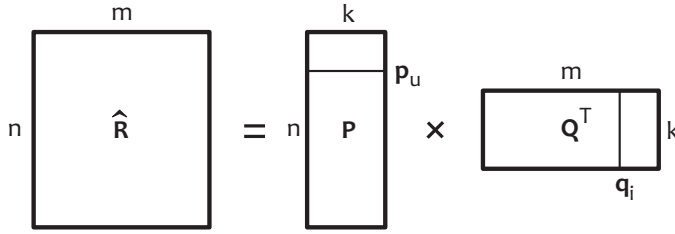


Рис. 5.13. Подход к совместной фильтрации на основе скрытых факторов

Если предположить, что число скрытых размерностей k фиксировано и $k \leq n$ и $k \leq m$, задача оптимизации 5.94 сводится к задаче низкоранговой аппроксимации, которую мы рассматривали в главе 2. Чтобы продемонстрировать подход к решению, допустим на минутку, что матрица рейтингов полная. В этом случае задача оптимизации имеет аналитическое решение в терминах сингулярного разложения (Singular Value Decomposition, SVD) матрицы рейтингов. В частности, с помощью стандартного алгоритма SVD матрицу можно разложить на произведение трех матриц:

$$\mathbf{R} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad (5.95)$$

где \mathbf{U} — матрица $n \times n$, ортонормированная по столбцам, $\mathbf{\Sigma}$ — диагональная матрица $n \times m$, а \mathbf{V} — матрица $m \times m$, ортонормированная по столбцам. Оптимальное решение задачи 5.94 можно получить в терминах этих факторов, усеченных до k наиболее значимых размерностей:

$$\hat{\mathbf{R}} = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T. \quad (5.96)$$

Следовательно, скрытые факторы, оптимальные с точки зрения точности прогнозирования, можно получить сингулярным разложением, как показано ниже:

$$\begin{aligned} \mathbf{P} &= \mathbf{U}_k \mathbf{\Sigma}_k, \\ \mathbf{Q} &= \mathbf{V}_k. \end{aligned} \quad (5.97)$$

Эта модель скрытых факторов, основанная на SVD, помогает решить проблемы совместной фильтрации, описанные в начале раздела. Во-первых, она заменяет большую матрицу рейтингов $n \times m$ матрицами факторов $n \times k$ и $m \times k$, которые обычно намного меньше, потому что на практике оптимальное количество скрытых размерностей k часто невелико. Например, известен случай, когда матрицу рейтингов с 500 000 пользователей и 17 000 элементов удалось достаточно хорошо аппроксимировать с использованием 40 измерений [Funk, 2016]. Далее, SVD устраняет корреляцию в матрице рейтингов: матрицы скрытых факторов, определяемые

выражением 5.97, являются ортонормированными по столбцам, то есть скрытые измерения не коррелированы. Если $k \ll n, m$, что обычно верно на практике, SVD также решает проблему разреженности, потому что сигнал, присутствующий в исходной матрице рейтингов, эффективно концентрируется (напомню, что мы выбираем k размерностей с наибольшей энергией сигнала), а матрицы скрытых факторов не разрежены. Рисунок 5.14 иллюстрирует это свойство. Алгоритм близости на основе пользователей (5.14, а) свертывает разреженные векторы рейтингов для данного элемента и данного пользователя, чтобы получить оценку рейтинга. Модель скрытых факторов (5.14, б), напротив, оценивает рейтинг путем свертки двух векторов уменьшенной размерности и с более высокой плотностью энергии.

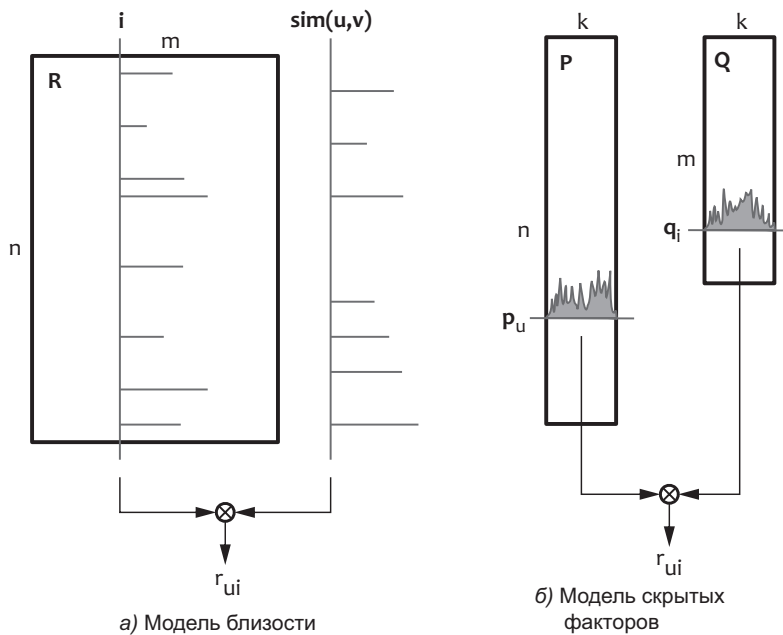


Рис. 5.14. Распределение энергии сигнала в моделях близости на основе пользователей и скрытых факторов

Только что описанный подход выглядит как стройное решение задачи скрытых факторов, но на самом деле он имеет серьезный недостаток из-за предположения о полноте матрицы рейтингов. Если матрица рейтингов разрежена, что имеет место почти всегда, стандартный алгоритм SVD нельзя применить напрямую, поскольку он не способен обрабатывать отсутствующие (неопределенные) элементы. Самым простым решением в этом случае является заполнение отсутствующих рейтингов некоторым значением по умолчанию, но это может привести к серьезному сме-

щению прогноза. Кроме того, это неэффективно с вычислительной точки зрения, потому что вычислительная сложность такого решения равна сложности SVD для полной матрицы $n \times m$, тогда как желательно иметь метод со сложностью, пропорциональной числу известных рейтингов. Эти проблемы можно решить с помощью альтернативных методов разложения, описанных в следующих разделах.

5.8.3.1. Разложение без ограничений

Стандартный алгоритм SVD — это аналитическое решение задачи низкоранговой аппроксимации. Однако эту проблему можно рассматривать как задачу оптимизации, и к ней также можно применить универсальные методы оптимизации. Один из самых простых подходов заключается в использовании метода градиентного спуска для итеративного уточнения значений скрытых факторов. Отправной точкой является определение функции стоимости J как остаточной ошибки прогноза:

$$\min_{P, Q} J = \| \mathbf{R} - \mathbf{PQ}^T \|^2. \quad (5.98)$$

Обратите внимание, что на этот раз мы не накладываем никаких ограничений, таких как ортогональность, на матрицы скрытых факторов. Вычисляя градиент функции стоимости по отношению к скрытым факторам, получаем следующий результат:

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{P}} &= -2(\mathbf{R} - \mathbf{PQ}^T)\mathbf{Q} = -2\mathbf{EQ}, \\ \frac{\partial J}{\partial \mathbf{Q}^T} &= -2\mathbf{P}^T(\mathbf{R} - \mathbf{PQ}^T) = -2\mathbf{P}^T\mathbf{E}, \end{aligned} \quad (5.99)$$

где \mathbf{E} — матрица остаточных ошибок:

$$\mathbf{E} = \mathbf{R} - \mathbf{PQ}^T. \quad (5.100)$$

Алгоритм градиентного спуска минимизирует функцию стоимости, перемещаясь на каждом шаге в отрицательном направлении градиента. Следовательно, можно найти скрытые факторы, минимизирующие квадрат ошибки прогнозирования рейтинга путем итеративного изменения матриц \mathbf{P} и \mathbf{Q} до сходимости, в соответствии со следующими выражениями:

$$\begin{aligned} \mathbf{P} &\leftarrow \mathbf{P} + \alpha \cdot \mathbf{EQ}, \\ \mathbf{Q}^T &\leftarrow \mathbf{Q}^T + \alpha \cdot \mathbf{P}^T\mathbf{E}, \end{aligned} \quad (5.101)$$

где α — *скорость обучения*. Недостатком метода градиентного спуска является необходимость вычисления всей матрицы остаточных ошибок и одновременного

изменения всех значений скрытых факторов в каждой итерации. Альтернативным подходом, который, возможно, лучше подходит для больших матриц, является стохастический градиентный спуск [Funk, 2016]. Алгоритм стохастического градиентного спуска использует тот факт, что общая ошибка прогноза J является суммой ошибок для отдельных элементов матрицы рейтингов, поэтому общий градиент J можно аппроксимировать градиентом в одной точке данных и изменять скрытые факторы поэлементно. Полная реализация этой идеи показана в алгоритме 5.1.

Алгоритм 5.1. Неограниченное разложение матрицы с применением алгоритма стохастического градиентного спуска. t — это счетчик итераций для среднего цикла, ε — порог сходимости, μ — среднее известных рейтингов, и $\mu_0 = \mu / \sqrt{k|\mu|}$

вход обучающее множество \mathbf{R} (выбранное из исходной матрицы рейтингов)

выход матрицы \mathbf{P} и \mathbf{Q}

инициализация $p_{ud}^{(0)} \sim$ случайное число со средним μ_0 $1 \leq u \leq m$, $1 \leq d \leq k$

инициализация $q_{id}^{(0)} \sim$ случайное число со средним μ_0 $1 \leq i \leq m$, $1 \leq d \leq k$

для размерности понятия $d=1, 2, \dots, k$ **выполнить**

повторять

для каждого рейтинга r_{ui} в обучающей выборке **выполнить**

$$\hat{r}_{ui} = \sum_{s=1}^d p_{us} \cdot q_{is}$$

$$e = r_{ui} - \hat{r}_{ui}$$

$$p_{ud} \leftarrow p_{ud} + \alpha \cdot e \cdot p_{id}$$

$$q_{id} \leftarrow q_{id} + \alpha \cdot e \cdot p_{ud}$$

$$SSE^{(t+1)} \leftarrow SSE^{(t)} + e^2 \text{ (сумма квадратов ошибок)}$$

конец

пока не выполнится условие $|SSE^{(t+1)} - SSE^{(t)}| < \varepsilon$ (условие сходимости)

конец

Первый этап алгоритма — инициализация матриц скрытых факторов. Выбор этих начальных значений не очень важен, но в данном случае выбрано равномерное распределение энергии известных рейтингов среди случайно сгенерированных скрытых факторов. Затем алгоритм последовательно оптимизирует размерности понятия. Для каждого измерения он многократно выполняет обход всех рейтингов

в обучающем наборе, прогнозирует каждый рейтинг с использованием текущих значений скрытых факторов, оценивает ошибку и корректирует значения факторов в соответствии с выражениями 5.101. Оптимизация измерения завершается по выполнении условия сходимости, после чего алгоритм переходит к следующему измерению.

Алгоритм 5.1 помогает преодолеть ограничения стандартного метода SVD. Он оптимизирует скрытые факторы, циклически перебирая отдельные точки данных, и тем самым избегает проблем с отсутствующими рейтингами и алгебраическими операциями с гигантскими матрицами. Итерационный подход также делает стохастический градиентный спуск более удобным для практических приложений, чем градиентный спуск, который изменяет целые матрицы с помощью выражений 5.101.

ПРИМЕР 5.6

По сути, подход на основе скрытых факторов — это целая группа методов обучения представлением, способных выявлять закономерности, неявно присутствующие в матрице рейтингов, и представлять их явно в виде понятий. Иногда понятия имеют вполне осмысленную интерпретацию, особенно высокоэнергетические, хотя это не означает, что все понятия всегда имеют осмысленное значение. Например, применение алгоритма разложения матриц к базе данных рейтингов фильмов может создать факторы, приблизительно соответствующие психографическим измерениям, таким как мелодрама, комедия, фильм ужасов и т. д. Проиллюстрируем это явление небольшим числовым примером, который использует матрицу рейтингов из табл. 5.3:

$$R = \begin{bmatrix} 5 & 4 & — & 1 & 2 & 1 \\ 4 & — & 3 & 1 & 1 & 2 \\ — & 5 & 5 & — & 3 & 3 \\ 2 & — & 1 & 4 & 5 & 4 \\ 2 & 2 & 2 & — & 4 & — \\ 1 & 2 & 1 & — & 5 & 4 \end{bmatrix}. \quad (5.102)$$

Сначала вычтем глобальное среднее $\mu = 2,82$ из всех элементов, чтобы центрировать матрицу, а затем выполним алгоритм 5.1 с $k = 3$ скрытыми измерениями и скоростью обучения $\alpha = 0,01$, чтобы получить следующие две матрицы факторов:

$$P = \begin{bmatrix} -1,40 & 0,30 & 0,95 \\ -1,03 & -0,90 & 0,34 \\ -0,94 & 1,53 & -0,12 \\ 1,26 & 0,34 & 0,66 \\ 0,80 & 0,16 & 0,07 \\ 1,47 & 0,38 & 0,05 \end{bmatrix} \quad Q = \begin{bmatrix} -1,16 & 0,31 & 0,60 \\ -0,96 & 0,82 & -0,43 \\ -1,26 & 0,71 & -0,83 \\ 1,29 & 0,30 & -0,37 \\ 1,18 & 0,90 & 0,60 \\ 0,83 & 0,37 & -0,44 \end{bmatrix}. \quad (5.103)$$

Каждая строка в этих матрицах соответствует пользователю или фильму, и все 12 векторов-строк изображены на рис. 5.15. Обратите внимание, что элементы в первом столбце (первый вектор понятий) имеют наибольшие величины, а величины в последующих столбцах постепенно уменьшаются. Это объясняется тем, что первый вектор-понятие захватывает столько энергии сигнала, сколько возможно захватить с помощью одного измерения, второй вектор-понятие захватывает только часть остаточной энергии и т. д. Далее, обратите внимание, что первое понятие можно семантически интерпретировать как ось драма — боевик, где положительное направление соответствует жанру боевика, а отрицательное — жанру драмы. Рейтинги в этом примере имеют высокую корреляцию, поэтому хорошо видно, что

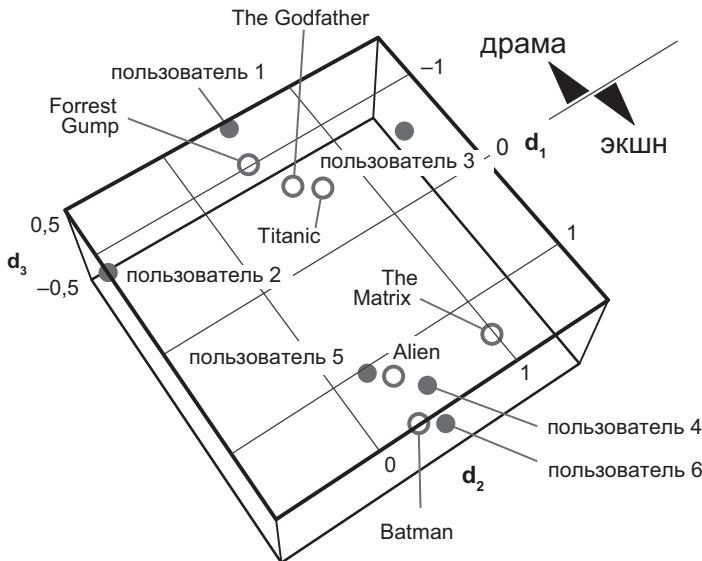


Рис. 5.15. Визуализация скрытых факторов 5.103. Измерения d_1 , d_2 и d_3 соответствуют первому, второму и третьему столбцам матриц соответственно. Пользователи и фильмы отображаются в виде заполненных и пустых кругов соответственно

первые три пользователя и первые три фильма имеют большие отрицательные значения в первом векторе-понятии (фильмы-драмы и пользователи, которым нравятся такие фильмы), тогда как три последних пользователя и три последних фильма имеют большие положительные значения в одном и том же столбце (боевики и пользователи, которые предпочитают этот жанр). Второе измерение в данном конкретном случае соответствует в основном смещению пользователя или элемента, которое можно интерпретировать как психографический атрибут (критичность суждений пользователя? популярность фильма?). Остальные понятия можно рассматривать как шум.

Полученные матрицы факторов не являются полностью ортогональными по столбцам, но стремятся к ортогональности, потому что это следует из оптимальности решения SVD. Это можно увидеть, рассматривая произведения $P^T P$ и $Q^T Q$, которые близки к диагональным матрицам:

$$\begin{aligned} P^T P &= \begin{bmatrix} \mathbf{8,28} & 0,19 & -0,62 \\ 0,19 & \mathbf{3,54} & 0,05 \\ -0,62 & 0,05 & \mathbf{1,47} \end{bmatrix}, \\ Q^T Q &= \begin{bmatrix} \mathbf{7,60} & -0,28 & 0,63 \\ -0,28 & \mathbf{2,31} & -0,49 \\ 0,63 & -0,49 & \mathbf{1,92} \end{bmatrix}. \end{aligned} \quad (5.104)$$

Матрицы 5.103 по сути являются предиктивной моделью, которую можно использовать для оценки как известных, так и отсутствующих рейтингов. Оценки можно получить путем умножения двух факторов и добавления обратно глобального среднего:

$$\begin{aligned} \hat{R} &= PQ^T + \mu = \\ &= \begin{bmatrix} 5,11 & 4,00 & [\mathbf{4,01}] & 0,75 & 2,00 & 1,35 \\ 3,94 & [\mathbf{2,93}] & 3,19 & 1,11 & 1,00 & 1,49 \\ [\mathbf{4,31}] & 5,03 & 5,19 & [\mathbf{2,12}] & 3,03 & 2,67 \\ 1,86 & [\mathbf{1,61}] & 0,94 & 4,30 & 5,01 & 3,71 \\ 1,99 & 2,15 & 1,88 & [\mathbf{3,87}] & 3,95 & [\mathbf{3,51}] \\ 1,26 & 1,69 & 1,20 & [\mathbf{4,81}] & 4,94 & 4,17 \end{bmatrix}. \end{aligned} \quad (5.105)$$

Результаты достаточно точно воспроизводят известные и предсказывают недостающие рейтинги в соответствии с интуитивными ожиданиями. Точность оценок можно увеличивать или уменьшать, изменяя число измерений, а оптимальное число измерений можно определить на практике путем перекрестной проверки и выбора разумного компромисса между вычислительной сложностью и точностью.

5.8.3.2. Разложение с ограничениями

Стандартный алгоритм SVD дает оптимальное решение задачи низкоранговой аппроксимации, а факторы \mathbf{P} и \mathbf{Q} , создаваемые им, являются ортогональными по столбцам. Алгоритм 5.1 стохастического градиентного спуска аппроксимирует это оптимальное решение. Если матрица рейтингов на входе полная, алгоритм 5.1 сходится к тем же результатам, ортогональным по столбцам, что и SVD. Единственное различие — диагональная масштабирующая матрица, присутствующая в SVD, представлена двумя факторами. Однако, если входная матрица рейтингов не является полной, выходные данные, полученные алгоритмом, могут быть не ортогональными. Это означает, что понятия остаются коррелированными в статистическом и геометрическом смыслах. Это может не влиять на качество прогнозирования рейтингов, но делает результаты менее интерпретируемыми из-за остаточных корреляций. Можно поставить вопрос: как наложить дополнительные ограничения, такие как ортогональность или неотрицательность, на скрытые факторы, чтобы лучше определить пространство понятий. К счастью, алгоритм стохастического градиентного спуска можно модифицировать для поддержки таких дополнительных ограничений и достижения существенной гибкости и контроля над процессом факторизации.

Рассмотрим задачу разложения с ограничениями ортогональности. Так же как в случае оптимизации без ограничений, цель состоит в том, чтобы найти скрытые факторы, минимизирующие квадрат ошибки прогноза, но при этом базис понятий должен быть ортогональным. Это приводит к следующей задаче оптимизации с ограничениями:

$$\begin{aligned} \min_{\mathbf{P}, \mathbf{Q}} \quad & \| \mathbf{R} - \mathbf{P} \cdot \mathbf{Q}^T \|^2 \\ \text{с учетом} \quad & \mathbf{P}^T \mathbf{P} \text{ является диагональной матрицей} \\ & \mathbf{Q}^T \mathbf{Q} \text{ является диагональной матрицей.} \end{aligned} \quad (5.106)$$

Метод градиентного спуска можно адаптировать к задачам оптимизации с ограничениями, применив метод, называемый проецируемым градиентным спуском. Идея метода заключается в применении ограничений на каждой итерации градиентного спуска, чтобы измененная переменная проецировалась обратно на множество допустимых решений. Более формально — градиентный спуск минимизирует функцию стоимости, итеративно перемещаясь в отрицательном направлении градиента:

$$\begin{aligned} \min_x \quad & f(x) \\ \text{правило обучения:} \quad & x^{(t+1)} = x^{(t)} - \alpha \cdot \nabla f(x^{(t)}). \end{aligned} \quad (5.107)$$

Проецируемый градиентный спуск распространяет этот подход на оптимизацию с ограничениями. На каждом шаге мы сначала перемещаемся в направлении отрицательного градиента, а затем корректируем решение, чтобы остаться в пределах допустимого множества:

$$\begin{aligned} \min_x f(x) \\ \text{с учетом } x \in C \\ \text{правило обучения: } z^{(t+1)} = x^{(t)} - \alpha \cdot \nabla f(x^{(t)}) \\ x^{(t+1)} = \operatorname{argmin}_{x \in C} \|z^{(t+1)} - x\|. \end{aligned} \quad (5.108)$$

В случае наложения ограничения ортогональности допустимым множеством для понятия d являются все векторы, ортогональные ранее вычисленным векторам понятий. Это означает, что решение, полученное градиентным спуском, можно отобразить на допустимое множество, вычитая его проекции на ранее вычисленные понятия. Например, если предположить, что первый вектор понятий p_1 пользователя (первый столбец матрицы P) определяется в ходе первой итерации внешнего цикла алгоритма 5.1, тогда решение-кандидат для второго вектора понятий p_2 (второй столбец матрицы P) можно ортогонализировать как

$$p_2 = p_2 - \operatorname{proj}(p_2, p_1), \quad (5.109)$$

где $\operatorname{proj}(a, b)$ — вектор проекции a на b , который определяется как

$$\operatorname{proj}(a, b) = \frac{a \cdot b}{b \cdot b} b. \quad (5.110)$$

На следующем шаге третий вектор понятий пользователя можно ортогонализировать путем вычитания его проекций на два предыдущих:

$$p_3 = p_3 - \operatorname{proj}(p_3, p_1) - \operatorname{proj}(p_3, p_2) \quad (5.111)$$

и так далее. Этот процесс, по сути, является итерационной версией процесса Грама–Шмидта (Gram–Schmidt), базовой процедуры в линейной алгебре, которая принимает произвольное множество линейно независимых векторов и строит на его основе множество ортогональных векторов. Тот же подход можно использовать для векторов понятий элемента, то есть для столбцов матрицы Q . Вставив эти операции ортогонализации в алгоритм 5.1, мы получим алгоритм 5.2, использующий точно такой же внутренний цикл для обновления элементов скрытых факторов, но имеющий дополнительный шаг проецирования для ортогонализации базиса векторов понятий.

Алгоритм 5.2. Разложение матрицы с ограничением ортогональности методом стохастического градиентного спуска

вход обучающее множество \mathbf{R} (выбранное из исходной матрицы рейтингов)

выход матрицы \mathbf{P} и \mathbf{Q}

инициализация $p_{ud}^{(0)} \sim$ случайное число со средним μ_0 $1 \leq u \leq n, 1 \leq d \leq k$

инициализация $q_{id}^{(0)} \sim$ случайное число со средним μ_0 $1 \leq i \leq m, 1 \leq d \leq k$

для размерности понятия $d = 1, 2, \dots, k$ **выполнить**

повторять

для каждого рейтинга r_{ui} в обучающей выборке **выполнить**

 | изменить элементы p_d и q_d (см. алгоритм 5.1)

конец

$$p_d \leftarrow p_d - \sum_{s=1}^{d-1} \text{proj}(p_d, p_s) \quad (\text{проекция})$$

$$q_d \leftarrow q_d - \sum_{s=1}^{d-1} \text{proj}(q_d, q_s) \quad (\text{проекция})$$

пока не выполнится условие сходимости

конец

Наконец, завершим этот раздел примечанием о разных типах ограничений, которые можно наложить на скрытые факторы, отличных от ортогональности. Алгоритм 5.2 создает строго ортогональные скрытые факторы даже в случае неполной матрицы рейтингов. Это до определенной степени улучшает интерпретируемость результатов, но отношения между пользователями, элементами и понятиями по-прежнему трудно интерпретировать из-за взаимовлияния положительных и отрицательных значений факторов. Можно попытаться решить эту проблему, заменив ограничение ортогональности ограничением неотрицательности:

$$\begin{aligned} \min_{\mathbf{P}, \mathbf{Q}} \quad & \|\mathbf{R} - \mathbf{P} \cdot \mathbf{Q}^T\|^2 \\ \text{с учетом} \quad & \mathbf{P} \geq 0 \\ & \mathbf{Q} \geq 0. \end{aligned} \tag{5.112}$$

Эту задачу оптимизации, известную как *неотрицательное матричное разложение*, также можно решить с помощью варианта алгоритма градиентного спуска [Zhang et al., 1996; Lee and Seung, 2001]. Преимуществом неотрицательного разложения

является лучшей интерпретируемость результатов, поскольку каждый элемент фактора указывает на близость к понятию, а каждого пользователя или элемент можно представить в виде аддитивной линейной комбинации понятий.

5.8.3.3. Продвинутое модели скрытых факторов

Методы разложения, рассмотренные в предыдущих разделах, обеспечивают прочную основу для создания моделей скрытых факторов. Однако эти алгоритмы очень просты и оставляют много места для улучшения и расширения. Такие расширенные модели, иногда очень сложные, могут обеспечить существенное улучшение качества основных k рекомендаций, даже притом что постепенное улучшение с точки зрения точности прогнозирования (например, среднеквадратичная ошибка) может быть очень ограниченным [Koren, 2008]. В этом разделе мы рассмотрим несколько продвинутых моделей, которые можно рассматривать как практические реализации подхода скрытых факторов. Эти методы базируются в основном на идеях, которые мы уже рассматривали в связи с другими алгоритмами рекомендаций, но были адаптированы к подходу скрытых факторов.

РЕГУЛЯРИЗАЦИЯ И СМЕЩЕНИЯ. Как говорилось в разделе 5.6.1, смещения пользователей и элементов являются важными *базовыми оценками*, способными выявить и удалить эффекты среднего пользователя и элемента. Поскольку базовые оценки и скрытые факторы можно определить с помощью градиентного спуска, мы можем объединить две модели в одну и совокупно оптимизировать смещения и переменные скрытых факторов. Формула прогнозирования рейтинга для этой модели определяется следующим образом:

$$\hat{r}_{ui} = \mu + b_i + b_u + \mathbf{p}_u \mathbf{q}_i^T, \quad (5.113)$$

где μ — глобальное среднее, b_i — смещение элемента, b_u — смещение пользователя, а последний член соответствует части модели, определяющей скрытые факторы. Добавляя регуляризационный член, который помогает избежать переобучения на разреженных данных, мы переводим эту модель в следующую задачу оптимизации:

$$\min \sum_{u,i} \left(r_{ui} - \mu - b_i - b_u - \mathbf{p}_u \mathbf{q}_i^T \right)^2 + \lambda \left(b_i^2 + b_u^2 + \|\mathbf{p}_u\|^2 + \|\mathbf{q}_i\|^2 \right), \quad (5.114)$$

где λ — параметр регуляризации, и минимизация выполняется одновременно по всем смещениям и переменным скрытых факторов. Эту задачу можно решить с помощью версии базового алгоритма 5.1 стохастического градиентного спуска со следующими правилами обучения:

$$\begin{aligned}
b_u &\leftarrow b_u + \alpha(e - \lambda \cdot b_u) \\
b_i &\leftarrow b_i + \alpha(e - \lambda \cdot b_i) \\
p_{ud} &\leftarrow p_{ud} + \alpha(e \cdot q_{id} - \lambda \cdot p_{ud}) \\
q_{id} &\leftarrow q_{id} + \alpha(e \cdot p_{ud} - \lambda \cdot q_{id}),
\end{aligned} \tag{5.115}$$

где α — скорость обучения. Эту модель можно рассматривать как практическую версию базового разложения без ограничений. Она часто упоминается как SVD-модель, что технически не совсем верно.

НЕЯВНАЯ ОБРАТНАЯ СВЯЗЬ. Вторая модель, которую мы рассмотрим, основывается на наблюдении, что пользователь выбирает элементы для оценки не случайно, а в соответствии с личными интересами и предпочтениями. Следовательно, полезный сигнал несут не только фактические значения рейтингов, но и позиции известных рейтингов (см. раздел 5.1.1). Мы можем выделить этот сигнал о взаимодействии пользователь — элемент из матрицы неявной обратной связи $n \times m$, которая содержит единицы в позициях известных оценок и нули в позициях отсутствующих оценок. Приводя каждую строку к единичной длине, мы определяем неявную матрицу обратной связи F как

$$f_{ui} = \begin{cases} |I_u|^{-1/2}, & \text{если } r_{ui} \text{ известно} \\ 0 & \text{в противном случае} \end{cases}, \tag{5.116}$$

где I_u — множество элементов, оцененных пользователем u . В общем случае матрица неявной обратной связи не обязательно является производной от матрицы рейтингов и может быть создана из другого источника данных. Например, матрицу неявной обратной связи можно создать на основе истории покупок или просмотра веб-страниц, в которой каждый ненулевой элемент будет указывать на взаимодействие между пользователем и элементом.

Идея модели разложения с неявной обратной связью состоит в том, чтобы ввести дополнительное множество факторов элементов, где каждое значение фактора y_{id} характеризует, *насколько акт оценки элемента i увеличивает или уменьшает близость к понятию d* . Обозначим это множество факторов матрицей Y с размерами $m \times k$. Произведение матрицы неявной обратной связи и этой новой матрицы элементов–факторов $FY = (z_{ud})$ является матрицей $n \times k$, в которой строки соответствуют пользователям, столбцы — понятиям, а каждый элемент z_{ud} можно интерпретировать как инкрементальную близость пользователя u к понятию d , согласно неявной обратной связи, то есть акту присваивания рейтинга. Эту инкрементальную близость можно непосредственно прибавить к основной матрице пользователи–факторы P , характеризующей близость поль-

зователь—понятие, полученную из значений рейтингов; результатом является задача оптимизации

$$\min_{\mathbf{P}, \mathbf{Q}, \mathbf{Y}} \|\mathbf{R} - (\mathbf{P} + \mathbf{F}\mathbf{Y})\mathbf{Q}^T\|^2. \quad (5.117)$$

Добавляя смещения пользователей и элементов, получаем следующую формулу прогнозирования рейтинга:

$$\hat{r}_{ui} = \mu + b_i + b_u + \left(\mathbf{p}_u + |\mathbf{I}_u|^{-1/2} \sum_{j \in \mathbf{I}_u} \mathbf{y}_j \right) \mathbf{q}_i^T, \quad (5.118)$$

где \mathbf{y}_j — строки матрицы \mathbf{Y} . Правила обучения для стохастического градиентного спуска можно прямо вывести из выражения 5.118, включая члены регуляризации и принимая градиенты:

$$\begin{aligned} b_u &\leftarrow b_u + \alpha(e - \lambda_1 \cdot b_u) \\ b_i &\leftarrow b_i + \alpha(e - \lambda_1 \cdot b_i) \\ p_{ud} &\leftarrow p_{ud} + \alpha(e \cdot q_{id} - \lambda_2 \cdot p_{ud}) \\ q_{id} &\leftarrow q_{id} + \alpha\left(e\left(p_{ud} + |\mathbf{I}_u|^{-1/2} \sum_{j \in \mathbf{I}_u} y_{jd}\right) - \lambda_2 \cdot q_{id}\right) \\ y_{jd} &\leftarrow y_{jd} + \alpha\left(e \cdot |\mathbf{I}_u|^{-1/2} \cdot q_{id} - \lambda_2 \cdot y_{jd}\right), \end{aligned} \quad (5.119)$$

где λ_1 и λ_2 являются параметрами регуляризации. Эта модель, известная как модель SVD++, может предложить более высокую точность, чем базовая модель SVD, из-за более точной обработки сигнала неявной обратной связи [Koren, 2008]. Модель SVD++ часто считается одной из самых продвинутых и эффективных моделей скрытых факторов.

СПЛАВ МЕТОДОВ СКРЫТЫХ ФАКТОРОВ С МЕТОДОМ БЛИЖАЙШИХ СОСЕДЕЙ. Наконец, рассмотрим модель, сочетающую разложение матриц с методом ближайших соседей. Как отмечалось в разделе 5.7.4, совместную фильтрацию на основе близости обычно можно рассматривать как задачу регрессии, которую можно решить и аналитически, и с помощью методов оптимизации, таких как стохастический градиентный спуск. Последний подход позволяет вернуть модель ближайших соседей в алгоритм разложения и оптимизировать скрытые факторы вместе с весами модели. Интегрированную модель можно получить, объединив выражения скрытых факторов 5.118 с выражением близости 5.69, которое мы рассмотрели ранее, и получить следующую формулу прогнозирования рейтингов:

$$\begin{aligned} \hat{r}_{ui} = & \mu + \underline{b}_u + \underline{b}_i + \left(\underline{p}_u + |I_u|^{-1/2} \sum_{j \in I_u} \underline{y}_j \right) \underline{q}_i^T + \\ & + |Q_{ui}^s|^{-1/2} \sum_{j \in Q_{ui}^s} \left((r_{uj} - b_{uj}) \underline{w}_{ij} + \underline{c}_{ij} \right), \end{aligned} \quad (5.120)$$

где Q_{ui}^s — соседи элемента i в множестве элементов, оцененных пользователем u (то есть s самых похожих элементов в I_u), а b_{ij} — базовый прогноз. Факторы $|I_u|^{-1/2}$ и $|Q_{ui}^s|^{-1/2}$ можно интерпретировать как надежность соответствующих членов, то есть количество рейтингов, на которых основана оценка, поэтому вклад термов масштабируется вверх или вниз соответственно. Затем ошибка прогнозирования рейтинга минимизируется по всем подчеркнутым переменным одновременно с помощью набора правил обучения, аналогичных набору 5.119, но с дополнительными правилами для весов w_{ij} и c_{ij} [Koren, 2008].

5.9. Гибридные методы

Выработка рекомендаций — обширная и сложная задача, поэтому идеальная рекомендательная система должна использовать несколько источников данных и учитывать широкий спектр эффектов и сигналов, таких как взаимодействие пользователей и элементов, сходство содержимого элементов и многие другие. Однако большинство методов рекомендаций могут использовать только один тип данных и фиксировать только определенный тип эффектов. Например, базовая совместная фильтрация ориентирована на анализ матрицы рейтингов и игнорирует содержимое элементов, тогда как фильтрация по содержимому действует с точностью до наоборот. То есть все методы имеют свои сильные и слабые стороны и могут дополнять друг друга. Гибридный подход пытается создать совершенные рекомендательные системы путем объединения нескольких базовых алгоритмов.

Мы уже видели несколько примеров объединения двух и более методов рекомендаций. Например, в разделе 5.7.4.3 методы ближайших соседей по пользователям и элементам объединялись с регрессионным анализом, а в разделе 5.8.3.3 мы дополнили базовую модель SVD данными неявной обратной связи. Такие гибридные решения, как правило, обеспечивают значительно более высокое качество, чем любой из составных алгоритмов. Наша следующая цель — разработать более систематическую и исчерпывающую основу для гибридизации, с использованием которой, в идеале, можно создавать оптимальные сочетания любых алгоритмов рекомендаций. Эта основа поможет нам не только создать более мощные службы рекомендаций, но и лучше понять некоторые из ранее описанных методов.

Проблема гибридных рекомендательных моделей тесно связана с ансамблевым обучением, сосредоточенным на методах, которые генерируют и объединяют несколько моделей классификации или регрессии, чтобы получить лучшее качество прогнозирования, чем позволяют отдельные алгоритмы обучения. В следующих разделах мы используем теорию ансамблей для создания гибридных моделей, начав с самых простых методов, и постепенно будем увеличивать сложность.

5.9.1. Переключение

Один из самых простых способов объединения нескольких рекомендательных алгоритмов — простое переключение между ними в зависимости от определенных условий. Например, можно предположить, что метод совместной фильтрации хорошо подходит для случаев, когда элемент имеет не слишком мало известных рейтингов, иначе с задачей прогнозирования лучше справится фильтрация по содержанию [Burke, 2002]. То есть мы можем организовать переключение между этими двумя алгоритмами в зависимости от количества пользователей, оценивших элемент:

$$\hat{r}_{ui} = \begin{cases} \hat{r}_{ui}^{(\text{collaborative})}, & \text{если } |U_i| > 20 \\ \hat{r}_{ui}^{(\text{content})} & \text{иначе} \end{cases}, \quad (5.121)$$

где U_i — множество пользователей, оценивших элемент i . Это решение может помочь обойти проблему холодного старта, характерную для совместной фильтрации, и в то же время улучшить тривиальные рекомендации, создаваемые фильтрацией по содержанию, когда это возможно. Общая схема такой модели переключения алгоритма рекомендаций показана на рис. 5.16. Данный подход, однако, несколько несовершенен, потому что основан на эвристических правилах, а не на формальном аппарате оптимизации. Мы определенно можем достичь лучших результатов, используя алгоритмы машинного обучения и оптимизации для правильного смешивания результатов отдельных моделей.

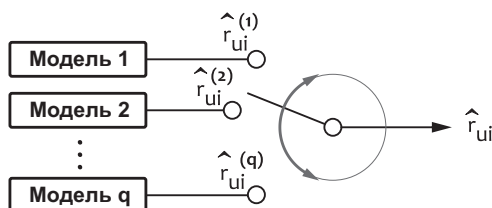


Рис. 5.16. Модель рекомендаций на основе переключения алгоритмов

5.9.2. Смешивание

Представьте несколько моделей рекомендаций, обученных на одном и том же наборе пользователей и элементов, то есть каждая из моделей может оценить рейтинг r_{ui} для данной пары пользователь/элемент. Наша задача — объединить эти оценки, чтобы получить окончательное значение рейтинга, в идеале более точное, чем прогнозы, сделанные какой-то одной из моделей. Это можно сделать с помощью эвристических правил, как в подходе переключения, описанном в предыдущем разделе, но в то же время это естественная задача регрессии, которую можно эффективно решить с помощью инструментов машинного обучения.

Задачу смешивания нескольких оценок рейтингов формально можно определить следующим образом. Предположим, что есть s обучающих образцов, то есть известных рейтингов в обучающем наборе. Этот набор используется для обучения q моделей рекомендаций, каждая из которых может предсказать рейтинг для данной пары пользователь/элемент. Для каждого обучающего образца j обозначим вектор результатов, возвращаемых моделью q (прогнозируемые значения рейтингов), как x_j , а истинное значение рейтинга — как y_j . Задачу смешивания доступных оценок можно определить как поиск функции смешивания $b(x)$, минимизирующей ошибку прогнозирования:

$$\min_b \sum_{j=1}^s (b(x_j) - y_j)^2. \quad (5.122)$$

Этот взгляд на задачу изображен на рис. 5.17. Задача объединения предсказаний нескольких алгоритмов обучения с использованием еще одного алгоритма обучения известна как *наложение* (stacking), поэтому мы будем использовать термины *смешивание* и *наложение* взаимозаменяемо. Наложение, по сути, является стандартной задачей обучения с учителем, которую можно решить с помощью разных алгоритмов классификации или регрессии.

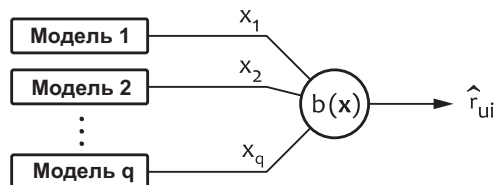


Рис 5.17. Смешивание алгоритмов рекомендаций

Одним из основных решений задачи 5.122, конечно, является линейная регрессия. В этом объединяющая функция является линейной и определяется как

$$b(\mathbf{x}) = \mathbf{x}^T \mathbf{w}, \quad (5.123)$$

где \mathbf{w} — вектор весов модели. Другими словами, конечное предсказанное значение рейтинга является линейной комбинацией предсказаний, сделанных отдельными алгоритмами рекомендаций:

$$\hat{r}_{ui} = \sum_{k=1}^q w_k \cdot \hat{r}_{ui}^{(k)}. \quad (5.124)$$

Оптимальные веса для функции смешивания можно рассчитать с использованием гребневой регрессии:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}, \quad (5.125)$$

где \mathbf{y} — вектор из s известных рейтингов, λ — параметр регуляризации, а \mathbf{X} — квадратная матрица прогнозов рейтингов. Каждый элемент x_{jk} — это рейтинг, предсказанный алгоритмом k для j -го обучающего образца.

На практике лучшие результаты часто дают нелинейные модели смешивания, такие как нейронные сети и деревья решений с градиентным бустингом [Jahger et al., 2010; Koren, 2009; Töschner et al., 2009]. Комбинация может включать десятки моделей рекомендаций разных типов (ближайших соседей, разложения матриц, смешанные и т. д.), и каждый тип может быть представлен несколькими вариантами модели, обученными с разными значениями числовых параметров, таких как число скрытых факторов. Модели могут обучаться на всем обучающем наборе, или этот набор может быть разделен на подмножества (блоки) случайным образом или в соответствии с некоторыми критериями, а затем отдельные модели могут быть обучены на каждом блоке отдельно. Смешивание является мощным методом, способным существенно улучшить качество прогнозов, поэтому многие методы смешивания были либо заимствованы из теории ансамблей, либо разработаны специально для рекомендательных систем. Мы рассмотрим некоторые из этих расширений и уточнений в следующих разделах.

5.9.2.1. Смешивание с последовательным обучением модели

Базовый подход к смешиванию предполагает, что все модели, участвующие в смешивании, обучаются заранее, а затем отдельно обучается функция смешивания, чтобы минимизировать общую ошибку прогнозирования. Однако такой подход не всегда является оптимальным, поскольку набор моделей, каждая из которых имеет низкую ошибку прогнозирования, не обязательно создаст сочетание с минимальной ошибкой прогнозирования. Отчасти это можно объяснить корреляцией между моделями — хорошее сочетание требует от моделей не только низких ошибок про-

гнозирования по отдельности, но и определенной степени независимости ошибок [Töscher et al., 2009]. Идеальным решением было бы объединить все модели в одну большую модель и одновременно оптимизировать все параметры относительно общей ошибки прогнозирования сочетания. Гибридная модель из раздела 5.8.3.3, объединяющая метод ближайших соседей со скрытыми факторами, на самом деле является примером такого решения. К сожалению, этот подход становится все более трудноразрешимым с увеличением числа моделей и, как следствие, количества параметров. С этой точки зрения смешивание можно рассматривать как аппроксимацию «разделяй и властвуй» глобального оптимального решения.

Смешивание можно улучшить, включив в процесс обучения модели функцию глобальной ошибки. Если предположить, что отдельные модели обучаются с использованием градиентного спуска, что верно для многих практических приложений, одним из возможных решений является переопределение условия сходимости цикла градиентного спуска как зависящего от общей ошибки прогнозирования. Данное решение реализовано в алгоритме 5.3.

Алгоритм 5.3. Последовательное обучение моделей с использованием линейной функции смешивания [Töscher et al., 2009]

$\mathbf{X}^{(0)} = s \times 1$ столбец матрицы из единиц (постоянный член)

для модели рекомендаций $k = 1, 2, \dots, q$ **выполнить**

повторять

 обновить модель k (один шаг обучения модели)

 предсказать рейтинги x_k с использованием модели k

$\mathbf{X} = [\mathbf{X}^{(k-1)} | \mathbf{x}_k]$ (попробовать добавить x_k в смесь)

$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ (оптимизировать функцию смешивания)

$\mathbf{r} = \mathbf{X} \cdot \mathbf{w}$ (вычислить предсказания для смеси)

$\text{SSE} = \|\mathbf{r} - \mathbf{y}\|^2$ (обновить общую ошибку предсказания)

пока не выполнится условие сходимости SSE

$\mathbf{X}^{(k)} = \mathbf{X}$ (постоянно добавлять x_k в смесь)

конец

Мы инициализируем матрицу \mathbf{X} выходов модели столбцом единиц, который можно интерпретировать как постоянный член. Затем модели рекомендаций обучаются одна за другой и добавляются в смесь. Каждая модель обучается с помощью градиентного спуска или стохастического градиентного спуска во внутреннем цикле

алгоритма 5.3. В каждой итерации мы обновляем модель, используя ее правила обучения, прогнозируем рейтинги для всех обучающих образцов, собираем временную матрицу \mathbf{X} , добавляя столбец с новыми прогнозируемыми рейтингами, повторно оптимизируем смесь (алгоритм использует функцию линейного смешивания лишь для иллюстрации, в действительности же может использоваться любая другая модель смешивания) и оцениваем ошибку прогнозирования всей смеси. Метод не меняет функций ошибок и правил обучения отдельных моделей, но меняет условие сходимости так, что обучение останавливается, когда общая ошибка прогнозирования смеси оказывается минимальной. На практике ошибка прогнозирования смеси может продолжать уменьшаться после того, как ошибка предсказания модели достигнет минимума и начнет увеличиваться.

5.9.2.2. Смешивание с остаточным обучением

С точки зрения независимости ошибок некоторые модели в смеси полезно обучать с использованием остаточных ошибок других моделей в роли входных данных. Для примера давайте объединим несколько моделей в цепочку так, чтобы результаты моделей в начале цепочки передавались на вход следующих за ними моделей, как показано на рис. 5.18.



Рис. 5.18. Обучение моделей на остаточных ошибках

Модели в цепочке обучаются последовательно. Первая модель обучается на исходных образцах. Рейтинги, предсказанные этой моделью, затем вычитаются из исходных образцов, и на этих остаточных ошибках обучается вторая модель, и т. д. Окончательная смесь может включать предсказания, созданные моделями, обученными на исходных данных и на остаточных ошибках. Среди рассмотренных ранее моделей примерами методов остаточного обучения могут служить удаление средних глобальных рейтингов и создание базовых предиктивных моделей.

5.9.2.3. Смешивание со взвешиванием признаков

Точность прогнозов рейтингов и, в конечном счете, качество рекомендаций можно повысить, объединив результаты нескольких моделей рекомендаций. В предыдущих разделах мы рассмотрели конструирование функции для оптимального смешивания результатов, возвращаемых моделями. Однако точность можно улучшить

еще больше, если функция смешивания использует не только результаты моделей рекомендаций, но и дополнительные сигналы о надежности модели или некоторые внешние сигналы о пользователях или элементах. Например, некоторые модели могут давать очень точные результаты, если для пользователя или элемента известно много рейтингов, но могут оказаться очень неточными и нестабильными при недостаточном количестве известных рейтингов. Вес таких моделей в смеси можно увеличивать или уменьшать в зависимости от распределения рейтингов. Некоторые следы этого подхода можно заметить в моделях, описанных выше. Например, метод взвешенного сходства, описанный в разделе 5.7.1, смешивает данные о надежности с мерами сходства пользователей. Аналогично модель скрытых факторов с неявной обратной связью, описанная в разделе 5.8.3.3, усиливает некоторые факторы, используя внешний неявный сигнал обратной связи.

В смешивании можно по-разному использовать преимущества внешних сигналов, иногда называемых *мета-признаками*. Например, сигналы можно передавать в функцию смешивания в виде дополнительных входных данных, то есть добавлять их в вектор результатов моделей рекомендаций. Хотя в целом этот подход вполне реализуем, он плохо работает с линейными функциями и часто требует обучения сложных нелинейных моделей смешивания, например с использованием деревьев решений с градиентным бустингом [Sill et al., 2009]. Альтернативное решение — объединить несколько линейных моделей в конвейер с предопределенной структурой, смешивающий сигналы от моделей рекомендаций с сигналами мета-признаков. Этот подход позволяет воспользоваться простотой и стабильностью линейной регрессии, но достичь гораздо лучших результатов, чем простая линейная модель, использующая мета-признаки в виде дополнительных входных данных. В остальной части раздела мы подробно обсудим детали этого метода [Sill et al., 2009].

Идея смешивания взвешиванием признаков заключается в смешивании результатов моделей рекомендаций с использованием линейной функции смешивания, вычисляющей весовые коэффициенты для смешивания на основе мета-признаков. Предположим, что множество из q моделей рекомендаций дает прогнозы рейтингов x_1, \dots, x_q . Допустим также, что каждому значению рейтинга соответствуют p мета-признаков g_1, \dots, g_p . Тогда смешать прогнозы можно с помощью функции линейного смешивания:

$$b(\mathbf{x}) = \sum_{k=1}^q w_k x_k, \quad (5.126)$$

с весовыми коэффициентами w_k , вычисляемыми динамически на основе мета-признаков:

$$w_k = f_k(\mathbf{g}) = \sum_{i=1}^p v_{ki} g_i, \quad (5.127)$$

где f_k — это то, что называют функциями признаков, а v_{ki} — статическими весами. Иначе говоря, функции признаков усиливают или подавляют сигналы от модели рекомендации, как показано на рис. 5.19. Обратите внимание, что эта конструкция очень похожа на конвейеры смешивания сигналов, рассмотренные в главе 4, в контексте служб поиска.

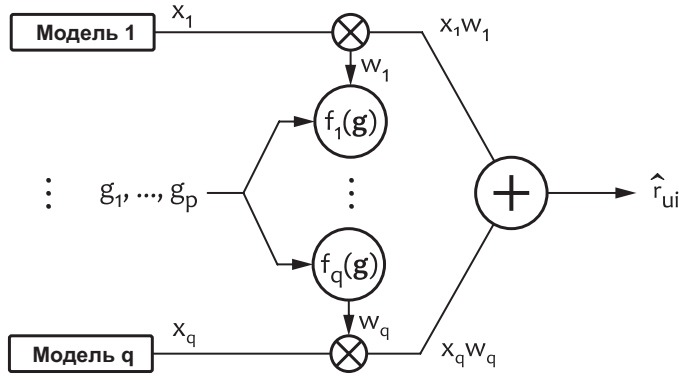


Рис. 5.19. Смешивание со взвешиванием признаков

Эта модель транслируется в следующую задачу оптимизации:

$$\min_v \sum_{j=1}^s \sum_{k,i} (v_{ki} \cdot g_{ji} \cdot x_{jk} - y_j)^2, \quad (5.128)$$

где внешняя сумма перебирает все s обучающих образцов, x_{jk} — рейтинг, предсказанный алгоритмом k для j -го обучающего образца, g_{ji} — мета-признаки для j -го образца, и y_j — истинное значение рейтинга для образца j . Чтобы решить эту задачу, введем $s \times (qp)$ матрицу A , получающуюся в результате векторного произведения предсказаний и мета-признаков для каждого обучающего образца:

$$\begin{aligned} a_{j,p(k-1)+i} &= x_{jk} \cdot g_{ji}, & 1 \leq j \leq s \\ & & 1 \leq k \leq q \\ & & 1 \leq i \leq p. \end{aligned} \quad (5.129)$$

В каждой строке матрицы A первые p элементов соответствуют первой модели рекомендаций, следующие p элементов соответствуют второй модели и т. д. Другими словами, первый сегмент p элементов является результатом первой модели x_{j1} , модулируемой каждой из функций признаков, и так далее. Введем также вектор-строку v с qp элементами, состоящий из весов v_{ki} и имеющий такую же

структуру, то есть первые p элементов соответствуют первой модели, следующие p элементов — второй модели и т. д. Оптимальные веса можно найти путем решения регрессионной задачи, соответствующей следующей системе линейных уравнений:

$$(A^T A + \lambda I) v = A^T y, \quad (5.130)$$

где λ — параметр регуляризации, а I — единичная диагональная матрица.

Смешивание со взвешиванием признаков — это относительно простое расширение базового взвешивания, позволяющее модулировать предсказания модели с помощью дополнительных сигналов или мета-признаков. Типичными примерами таких мета-признаков могут служить базовые статистики (например, сколько раз элемент был оценен пользователями, стандартное отклонение рейтингов пользователей), зависящие от времени статистики (например, количество разных дат, когда пользователь присваивал рейтинги элементам) и статистика корреляции (например, максимальная корреляция элемента с любым другим элементом). Эти статистики имеют большое значение для гибридной модели рекомендаций, потому что надежность оценок, генерируемых составными моделями, зависит от числа имеющихся рейтингов и других аналогичных факторов, отраженных в статистических сигналах. Следовательно, эти признаки позволяют гибридной модели научиться переключаться между моделями в зависимости от ожидаемой надежности оценок.

5.9.3. Расширение признаков

Следующий класс гибридных методов, который мы рассмотрим, — методы рекомендаций с расширением признаков. Метод расширения признаков относится к архитектуре с несколькими моделями рекомендаций, объединенными в цепочку так, что прогнозы, полученные одной моделью рекомендаций, используются другой моделью в качестве входных данных. Мы уже использовали этот подход в сочетании с остаточным обучением, но эту идею можно также использовать немного иначе.

Один из возможных подходов к объединению двух моделей рекомендаций в цепочку — использование первой модели в цепочке для создания совершенно новых признаков, которых нет в исходных данных, чтобы следующие модели могли использовать их в качестве входных данных. Например, наивный байесовский классификатор по содержанию, описанный в разделе 5.5.2, для рекомендации книг может использовать атрибуты элемента, такие как *похожие авторы* и *похожие названия*. Эти атрибуты можно создать с использованием меры сходства элементов, вычисленной по матрице рейтингов, то есть с помощью совместной

фильтрации [Mooney and Roy, 1999]. В результате получается гибридная модель с расширением признаков, где совместная фильтрация является первой моделью в цепочке, которая генерирует новые признаки, а классификатор по содержанию является второй моделью, которая использует эти признаки.

Второй вариант объединения моделей в цепочку — использование первой модели в цепочке для улучшения признаков, присутствующих в исходных данных. Например, для заполнения отсутствующих элементов в матрице рейтингов можно использовать модель рекомендаций по содержанию, а затем эту улучшенную матрицу можно передать какому-либо методу совместной фильтрации. Сравните этот подход с предыдущим примером, где совместная фильтрация использовалась с целью расширения входных данных для наивного байесовского классификатора по содержанию. Давайте конкретизируем это решение, называемое *совместной фильтрацией с бустингом по содержанию*, предположив, что в качестве компонента совместной фильтрации в гибриде используется модель ближайших соседей на основе пользователей [Melville et al., 2002]. Первый шаг — использование модели фильтрации по содержанию для заполнения отсутствующих элементов матрицы рейтингов и создания новой матрицы псевдорейтингов:

$$z_{ui} = \begin{cases} r_{ui}, & \text{если пользователь } u \text{ оценил элемент } i \\ \hat{r}_{ui}^{(c)} & \text{в противном случае,} \end{cases} \quad (5.131)$$

где $\hat{r}_{ui}^{(c)}$ — рейтинг, прогнозируемый моделью рекомендаций по содержанию. Второй шаг — применение компонента гибрида, реализующего совместную фильтрацию, к матрице псевдорейтингов с целью прогнозирования рейтинга для данной пары пользователь/элемент. В принципе, для окончательного прогноза можно использовать любой готовый алгоритм совместной фильтрации. Однако проблема в том, что рейтинги, добавленные на предыдущем шаге, искажают статистику количества известных рейтингов, используемых многими алгоритмами совместной фильтрации. Это требует изменить часть, отвечающую за совместную фильтрацию, и ввести несколько дополнительных факторов и параметров, чтобы исправить статистику:

- Надежность прогнозов рейтингов на основе содержимого зависит от количества известных рейтингов для данного пользователя. Следовательно, прогнозы, не имеющие достаточной поддержки, следует девальвировать при использовании совместной фильтрации. Если на шаге совместной фильтрации используется модель на основе близости пользователей, мы можем учесть надежность входящих рейтингов, изменив меру сходства пользователей. Давайте сначала определим нормализованную переменную поддержки, растущую пропорционально количеству рейтингов, присвоенных пользователем,

но ограниченную единицей, если количество оценок превышает пороговый параметр T :

$$q_u = \begin{cases} 1, & |I_u| \geq T \\ |I_u|/T, & \text{иначе} \end{cases}. \quad (5.132)$$

Затем переопределим функцию сходства, добавив коэффициент, равный гармоническому среднему из переменных поддержки для двух пользователей:

$$\text{sim}'(u, v) = \frac{2q_u q_v}{q_u + q_v} \cdot \text{sim}(u, v). \quad (5.133)$$

Гармоническое среднее выбрано потому, что оно смещено к минимуму двух чисел, благодаря этому мера сходства будет значительно уменьшена, если любой из пользователей присвоит слишком мало рейтингов.

- Гибридная система включает компоненты фильтрации по содержимому и совместной фильтрации, которые способны предсказывать рейтинги для данного пользователя и элемента, поэтому эти два прогноза должны смешиваться. Для управления балансом между прогнозами введем коэффициент усиления w_u прогноза по содержимому. Этот коэффициент определяется как базовый вес усиления w_{\max} , умноженный на переменную поддержки, чтобы уменьшить влияние ненадежных прогнозов:

$$w_u = w_{\max} \cdot q_u. \quad (5.134)$$

Окончательную формулу прогнозирования рейтинга для части, выполняющей совместную фильтрацию, можно определить следующим образом:

$$\hat{r}_{ui} = \mu_u + \frac{w_u (\hat{r}_{ui}^{(c)} - \mu_u) + \sum_v \text{sim}'(u, v) (z_{vi} - \mu_v)}{w_u + \sum_v \text{sim}'(u, v)}, \quad (5.135)$$

где μ_u — средний рейтинг пользователя, вычисленный по матрице псевдорейтингов. По сути это базовая модель ближайшего соседа на основе пользователей с добавлением оценки по содержимому $\hat{r}_{ui}^{(c)}$ и функцией сходства, скорректированной с учетом надежности этих оценок. Обратите внимание, что эти корректировки, связанные с надежностью, очень похожи на смешивание со взвешиванием признаков, рассматривавшихся в предыдущем разделе: в обоих случаях гибридная модель рекомендаций использует статистики рейтинга для девальвации моделей с низкой надежностью и усиления сигналов от моделей с высокой надежностью.

Модель совместной фильтрации по содержанию, определяемая выражением 5.135, как правило, дает более высокую точность, чем любой из двух составляющих ее методов. Если исходная матрица рейтингов достаточно плотная, модель превосходит оба компонента — фильтрации по содержанию и совместной фильтрации, — используя преимущество двух сигналов. В случае с разреженной матрицей рейтингов точность компонента совместной фильтрации снижается, а общее качество гибридной модели сходится к точности метода рекомендаций по содержанию [Melville et al., 2002].

5.9.4. Варианты представления гибридных рекомендаций

Завершим обзор гибридных методов краткой ремаркой о том, как гибридная модель рекомендаций может использовать презентационные возможности службы рекомендаций. Во-первых, стоит отметить, что рекомендации, подготовленные различными моделями, не обязательно смешивать вместе — система рекомендаций может просто отобразить несколько списков рекомендуемых элементов. Например, на веб-сайтах электронной коммерции часто отображается сразу несколько панелей рекомендаций с разными семантическими значениями, такие как *Вместе с этим элементом другие клиенты также просматривали*, *По вашей истории просмотров*, *Высоко оцененные элементы*, *Похожие элементы* и т. д. Естественно, эти панели можно заполнять с использованием разных алгоритмов рекомендаций, включая персонализированные и неперсонализированные. Системы рекомендаций, использующие этот подход, обычно называются *смешанными гибридами*.

В некоторых случаях рекомендуемые элементы должны отвечать дополнительным требованиям или условиям, в зависимости от того, как представлены рекомендации и как пользователь взаимодействует со службой рекомендаций. Например, пользователь может явно запросить рекомендацию ресторанов в определенном месте или книг, похожих на выбранную. В таких случаях рекомендательная система может использоваться как компонент, обрабатывающий результаты, полученные службой поиска или другой моделью рекомендаций. Например, служба поиска может использоваться для получения списка элементов, соответствующих критериям пользователя, а затем этот список можно отсортировать с помощью модуля совместной фильтрации. Этот метод, иногда называемый *каскадным*, можно рассматривать как крайний случай смешивания, когда сигнал от первой модели или службы поиска используется для разбивки элементов на две группы — релевантных и нерелевантных — элементов, а вторая модель рекомендаций выполняет вторичную сортировку внутри групп.

5.10. Контекстные рекомендации

Большинство рекомендательных алгоритмов, включая все методы, рассмотренные выше в этой главе, основаны на предположении, что релевантность данного элемента для данного пользователя можно предсказать, используя только профили элемента и пользователя. При таком подходе полностью игнорируются ситуации, когда рекомендации должны даваться с учетом времени, местоположения пользователя, маркетингового канала и другой информации о ситуации и окружающей среде. Однако эта контекстная информация очень важна, потому что потребители почти всегда принимают решения, исходя из контекста. Следовательно, релевантность рекомендаций уникальна для каждого отдельного случая рекомендаций и не определяется статистическими характеристиками элементов и профилей пользователей. Рассмотрим несколько конкретных случаев, чтобы лучше понять понятие контекста.

МЕСТОПОЛОЖЕНИЕ. Рекомендации обувного магазина для пользователей из Мурманска могут оказаться не актуальными для пользователей из Сочи. Клиенты, использующие мобильное приложение для поиска ресторанов поблизости, могут получить нерелевантные рекомендации после переезда в другое место.

ВРЕМЯ. Рекомендации фильмов, релевантные для пятнадцатилетнего пользователя сегодня, могут стать нерелевантными через пять лет, когда этому пользователю будет двадцать. Рекомендации телевизионных программ, релевантные утром, могут стать нерелевантными вечером. Рекомендации, сделанные системой рекомендаций в один сезон, могут потерять релевантность в другом сезоне.

НАМЕРЕНИЕ. Релевантность рекомендаций ресторана может меняться в зависимости от того, с кем пользователь ужинает: один, с супругом, с коллегами или с семьей. Рекомендации, сгенерированные для пользователей, совершающих покупки для себя, могут отличаться от рекомендаций, когда пользователю требуется купить подарок. Рекомендации по бронированию гостиниц для командированных могут отличаться от рекомендаций для туристов.

КАНАЛ. Рекомендации в электронных письмах могут иметь структуру и представление, отличные от структуры и представления рекомендаций на веб-сайте или в магазине.

УСЛОВИЯ. Рекомендации, сделанные универмагом, могут включать или не включать зонтики, в зависимости от текущих или прогнозируемых погодных условий.

Рекомендательная система должна принимать во внимание контекстную информацию о местоположении, времени, намерении и канале, чтобы обеспечить реле-

вантность в реальном времени. В оставшейся части этого раздела мы обсудим, как можно расширить или модифицировать алгоритмы рекомендаций, чтобы включить в них эти контекстные сигналы.

5.10.1. Многомерная основа

Традиционные модели рекомендаций предсказывают релевантность данного элемента для данного пользователя, опираясь на профили этого элемента и пользователя. Эти модели можно рассматривать как функции, получающие пользователя и элемент в качестве аргументов и прогнозирующие значение рейтинга:

$$\hat{r}_{ui} = R(u, i). \quad (5.136)$$

Контекстно-зависимые рекомендательные системы расширяют эту основу дополнительными аргументами, каждый из которых представляет определенное измерение контекста, такое как местоположение, время или канал [Adomavicius and Tuzhilin, 2008]:

$$\hat{r}_{ui} = R(u, i, \text{location}, \text{time}, \dots). \quad (5.137)$$

Другими словами, базовая функция рейтинга, определенная на двухмерном пространстве

$$R: \text{User} \times \text{Item} \rightarrow \text{Rating} \quad (5.138)$$

заменяется функцией, определенной на многомерном пространстве, включающем измерения «пользователь», «элемент», а также измерения, определяющие контекст:

$$R: \text{User} \times \text{Item} \times \text{Location} \times \text{Time} \times \dots \rightarrow \text{Rating}. \quad (5.139)$$

Эту идею иллюстрирует рис. 5.20, где изображен пример трехмерного пространства рекомендаций. В этом примере каждое значение рейтинга является функцией от пользователя, элемента и времени. Все известные рейтинги приписываются временной метке и помещаются в соответствующие ячейки трехмерного массива, а не в двухмерную матрицу рейтингов. Таким образом, целью модели рекомендаций является прогнозирование значений рейтинга в пустых ячейках массива. Обратите внимание, что многомерный массив можно свернуть в стандартную двухмерную матрицу рейтингов, отбросив контекстную информацию. Для этого может потребоваться объединение нескольких значений рейтинга, проецируемых в один элемент матрицы. Например, если пользователь оценил один и тот же товар несколько раз в разные даты, в матрице рейтингов можно сохранить только последнее или среднее

значение. Кроме того, матрицу рейтингов можно получить, выбрав определенную точку в измерении контекста и получив для этой точки двухмерный срез многомерного куба. Например, массив, изображенный на рис. 5.20, можно рассматривать как стопку матриц рейтингов $R(t)$, по одной для каждого временного интервала. Наконец, матрицу рейтингов можно создать не для определенной точки, а для определенного диапазона в контекстном измерении. В случае с примером на рис. 5.20 это мог бы быть некоторый интервал времени.

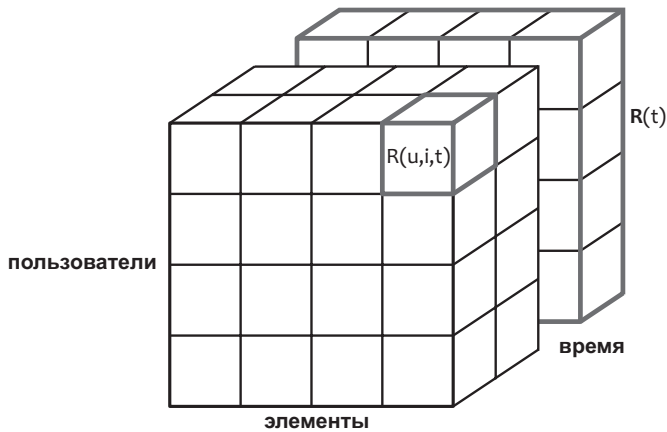


Рис. 5.20. Пример трехмерного пространства рекомендаций

Обычно предполагается, что контекстные измерения могут иметь иерархическую структуру. Например, каждый известный рейтинг может быть связан с датой, когда он был присвоен, поэтому измерение времени дискретно и содержит столько интервалов, сколько существует различных меток даты в данных рейтингов. Даты, однако, можно агрегировать в еженедельные, ежемесячные, квартальные или годовые интервалы, и, соответственно, матрицу рейтингов $R(t)$ можно сократить до определенной недели, месяца, квартала или года. Для одного измерения может быть несколько иерархий. Например, даты можно разделить на будние и выходные дни, чтобы получить матрицы рейтингов для будних и выходных дней. Аналогично можно объединить подробные атрибуты местоположения, такие как широта и долгота, в почтовые индексы, города, штаты и страны. Наконец, измерения «пользователь» и «элемент» также могут иметь определенную иерархию. Например, пользователей можно разделить на группы по возрасту, а элементы — по классам или видам.

Контекстную информацию можно получить из разных источников. Некоторые атрибуты, такие как время присваивания рейтинга или местоположение устройства пользователя, могут автоматически извлекаться системой рекомендаций или мар-

кетинговыми каналами, с которыми связана система. Некоторые другие атрибуты, особенно имеющие отношение к намерениям, часто недоступны непосредственно, но могут определяться с помощью специальных признаков в пользовательском интерфейсе (например, флажок *Подарочный заказ* в форме онлайн-заказа) или прогнозирующих моделей.

5.10.2. Контекстно-зависимые методы рекомендаций

Службу рекомендаций, не зависящую от контекста, можно рассматривать как процесс, который использует обучающие данные в форме $User \times Item \times Rating$, создает модель, отображающую пары пользователь/элемент в рейтинг, и использует эту модель для данного пользователя, чтобы создать отсортированный список рекомендаций. Этот конвейер показан на рис. 5.21.

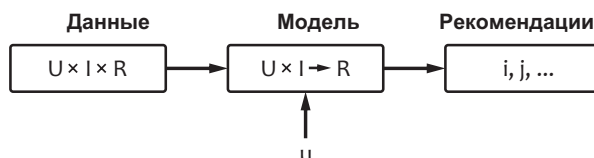


Рис. 5.21. Основные этапы процесса рекомендаций, не зависящего от контекста [Adomavicius and Tuzhilin, 2008]. U , I и R — измерения «пользователи», «элементы» и «рейтинги» соответственно. Рекомендуемые элементы обозначаются как i, j

Многомерная основа, описанная в предыдущем разделе, предлагает несколько идей, как можно модифицировать этот конвейер для включения контекстной информации [Adomavicius and Tuzhilin, 2008].

ПРЕДВАРИТЕЛЬНАЯ КОНТЕКСТНАЯ ФИЛЬТРАЦИЯ. Первое возможное решение — создание двухмерной матрицы рейтингов из исходных многомерных данных с последующим применением стандартной модели рекомендаций без учета контекста или алгоритма, принимающего эту матрицу на входе, как показано на рис. 5.22. В данном случае матрица рейтингов — это срез многомерного куба для заданного значения контекста. Например, службу рекомендаций фильмов, хранящую рейтинги с метками времени, можно превратить в службу рекомендаций для будних и выходных дней, используя две разные матрицы. Матрицы создаются путем выбора из исходного куба данных всех рейтингов, присвоенных в будние или в выходные дни соответственно.

Одной из ключевых проблем в методе предварительной фильтрации является поиск компромисса между разреженностью данных и точностью согласования с контекстом. С одной стороны, согласование входных данных с контекстом

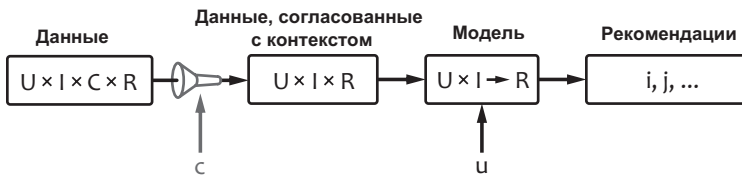


Рис. 5.22. Контекстно-зависимая система рекомендаций с предварительной контекстной фильтрацией [Adomavicius and Tuzhilin, 2008]

может повысить точность рекомендаций, поскольку алгоритм использует только рейтинги, релевантные контексту. С другой стороны, такое согласование уменьшает объем данных, доступных для алгоритма, что может отрицательно сказаться на качестве рекомендаций. Например, служба рекомендаций фильмов, которая дает рекомендации для просмотра фильмов в выходные дни, используя только рейтинги, присвоенные в выходные, почти наверняка потеряет некоторые релевантные сигналы, передаваемые рейтингами будних дней. Согласование с узкими критериями отбора также может привести к получению очень разреженных матриц рейтингов и, соответственно, к менее надежным и искаженным прогнозам рейтингов. Компромисс между разреженностью данных и точностью контекста можно контролировать с помощью иерархического агрегирования, о котором говорилось в предыдущем разделе. Например, зависимость от времени можно применить на уровне детализации дней недели (всего семь сегментов) или на уровне выходных дней (всего два сегмента). Обратите внимание, что, в принципе, можно использовать метод контролируемого снижения точности, обсуждавшийся в разделе 4.6.2, в контексте служб поиска, чтобы попробовать разные уровни детализации и выбрать оптимальный для данного значения контекста.

ЗАКЛЮЧИТЕЛЬНАЯ КОНТЕКСТНАЯ ФИЛЬТРАЦИЯ. Альтернативный подход к согласованию с контекстом — заключительная фильтрация рекомендаций, как показано на рис. 5.23. Система рекомендаций с заключительной фильтрацией первоначально не учитывает контекстную информацию и сворачивает куб данных с рейтингами в плоскую матрицу, к которой затем применяет стандартный алгоритм для получения рекомендаций без учета контекста. Затем список рекомендуемых элементов согласуется с контекстом с помощью контекстно-зависимых правил постобработки. Эти правила обычно основаны на атрибутах пользователя или элемента и могут быть эвристическими или управляться предиктивной моделью. Например, система рекомендаций одежды может создать начальный список несогласованных рекомендаций с помощью любого алгоритма — фильтрации по содержанию или совместной фильтрации, — а затем отфильтровать или изменить

ранг элементов в соответствии с текущим сезоном или погодой. Например, зимой она может поднимать вверх теплую одежду. В этом случае теплую одежду можно определить с помощью эвристического правила или модели классификации содержимого, например наивного байесовского классификатора.

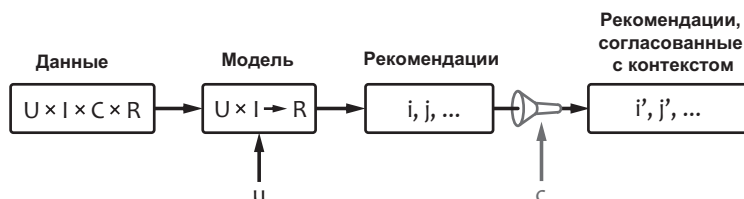


Рис. 5.23. Контекстно-зависимая система рекомендаций с заключительной контекстной фильтрацией [Adomavicius and Tuzhilin, 2008]

КОНТЕКСТНОЕ МОДЕЛИРОВАНИЕ. Наиболее общим решением задачи контекстных рекомендаций является создание модели, способной предсказывать рейтинги непосредственно как функции нескольких аргументов, включая элемент, пользователя и контекст. Такой подход показан на рис. 5.24. Ключевым преимуществом контекстного моделирования является возможность изучения и оптимизации контекстных параметров вместе с другими частями модели. Это может привести к лучшим результатам, чем могут дать эвристические решения с предварительной и заключительной фильтрацией, которые, как говорилось выше, могут отфильтровать релевантные данные и сигналы, формально не соответствующие критериям контекста.

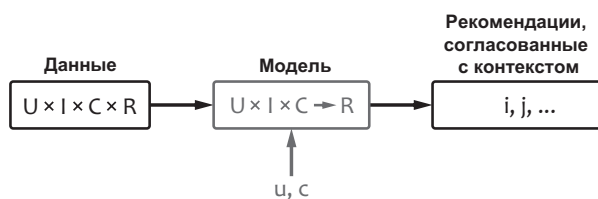


Рис. 5.24. Контекстно-зависимая система рекомендаций с контекстным моделированием [Adomavicius and Tuzhilin, 2008]

Контекстные модели часто можно получить расширением стандартной фильтрации по содержимому или совместной фильтрации. Проиллюстрируем эту идею на концептуальном примере использования модели ближайших соседей. Напомню, что модель ближайших соседей, модель, не согласованную с контекстом, можно выразить следующим образом:

$$\hat{r}_{ui} = \sum_{v,j} \text{sim}((u, i), (v, j)) \cdot r_{vj}, \quad (5.140)$$

где v и j — индексы ближайших пользователей и элементов соответственно. Мету сходства между парами пользователь/элемент (u, i) и (v, j) можно задать по-разному, как показывалось выше, в разделе, посвященном совместной фильтрации на основе близости. В случае модели ближайших соседей по пользователям, например, мерой будет служить сходство между двумя пользователями, которое, в свою очередь, можно вычислить как коэффициент корреляции Пирсона или как-то иначе:

$$\text{sim}((u, i), (v, j)) = \begin{cases} \text{sim}(u, v), & i = j \\ 0 & \text{в противном случае.} \end{cases} \quad (5.141)$$

Контекстная модель, учитывающая время, может расширить понятие сходства и включить измерение времени:

$$\hat{r}_{uit} = \sum_{v,j,s} \text{sim}((u, i, t), (v, j, s)) \cdot r_{vjs}, \quad (5.142)$$

где s — индекс рейтингов с метками времени, близкими к целевой контекстной метке времени t . Многомерную мету сходства можно определить как расстояние между двумя ячейками в трехмерном кубе данных. Например, можно использовать метрику евклидова расстояния:

$$\begin{aligned} \text{sim}((u, i, t), (v, j, s)) &= \\ &= \sqrt{\text{sim}^2(u, v) + \text{sim}^2(i, j) + \text{sim}^2(t, s)}. \end{aligned} \quad (5.143)$$

Это был всего лишь концептуальный пример, который иллюстрирует подход к включению контекстной информации в модель рекомендаций, но в следующем разделе мы разработаем более практические решения.

Наконец, следует отметить, что в гибридную модель можно объединить несколько алгоритмов рекомендаций, зависимых и не зависимых от контекста. Гибридный подход может помочь преодолеть ограничения отдельных методов (например, проблемы разреженности, вызванные чрезмерно строгим согласованием входных данных с контекстом) путем оптимального смешивания нескольких прогнозов.

5.10.3. Модели рекомендаций с учетом времени

Размерность, связанная с временем, является одним из наиболее важных типов контекста из-за существенной изменчивости шаблонов взаимоотношений между пользователями и элементами с течением времени. Хорошим примером может

служить изменение популярности элементов с течением времени, обусловленное внешними факторами, такими как изменения в моде, делающие некоторые предметы одежды более или менее популярными, участие актера в новом фильме, увеличивающее популярность похожих фильмов, или новые достижения в технологиях, делающие электронные устройства устаревшими. Другим примером может служить изменение пользовательских предпочтений со временем, которое может быть вызвано изменением вкусов, социальной роли или места жительства. Например, пользователь, прежде присваивавший ничем не выделяющимся элементам 4 звезды, со временем может стать более критичным в своих оценках и присваивать тем же элементам только 3 звезды. В то же время информация о временном контексте также является одним из самых простых для сбора типов данных, потому что временные метки могут устанавливаться внутри рекомендательной системы, без внешних зависимостей от маркетинговых каналов или пользовательских интерфейсов. Эти факторы делают рекомендации с учетом времени легко реализуемыми и могут существенно повысить точность рейтинговых прогнозов при относительно низких затратах.

Многомерную основу и методы согласования с контекстом на основе предварительной и заключительной фильтрации с легкостью можно применить к случаям с периодическим временным контекстом, таким как сезонность. Однако эти инструменты слишком просты, поэтому далее мы посмотрим, как воспользоваться преимуществами методов прогнозирования и оптимизации для создания более продвинутых и точных моделей рекомендаций с учетом времени. Мы увидим, как три основные модели совместной фильтрации: базовых оценок, ближайших соседей и скрытых факторов — можно расширить, чтобы учесть в них измерение времени. Стоит отметить, что все три решения были разработаны как компоненты единой гибридной модели [Koren, 2009]. Однако каждая модель использует собственный метод учета изменений во времени.

5.10.3.1. Базовые оценки с учетом временной динамики

Цель базовых оценок — зафиксировать закономерности, свойственные среднему пользователю, и смещения в рейтингах элементов, а также средний глобальный рейтинг. Напомним, что стандартная базовая оценка рейтинга определяется следующим образом (раздел 5.6.1):

$$b_{ui} = \mu + b_u + b_i, \quad (5.144)$$

где b_u и b_i — смещения среднего пользователя и элемента соответственно. Предположив дополнительно, что смещения пользователя и элемента могут изменяться с течением времени, версию с учетом времени можно определить как

$$b_{ui} = \mu + b_u(t) + b_i(t), \quad (5.145)$$

где $b_u(t)$ и $b_i(t)$ — функции времени, которые необходимо определить по имеющимся данным. Параметр t можно задать как количество дней, прошедших от некоторой нулевой даты в прошлом. На практике компоненты, представляющие пользователей и элементы, могут иметь очень разную временную динамику и свойства, поэтому для этих двух функций могут потребоваться два разных решения [Коген, 2009]. На практике популярность элементов медленно меняется с течением времени, и каждый элемент имеет относительно много рейтингов. Соответственно, временной диапазон можно разделить на несколько временных интервалов (например, по несколько недель в каждом) и для каждого независимо оценить смещение элемента. В результате мы приходим к следующей простой модели учета для элемента:

$$b_i = b_i + b_{i,\Delta t}, \quad (5.146)$$

где b_i — глобальное стационарное смещение, Δt — временной интервал, в который попадает t , и $b_{i,\Delta t}$ — смещение оценки элемента в течение этого периода. Обратите внимание, что часть смещения, зависящую от времени, можно оценить только для дат в прошлом и ей следует присвоить ноль, если весь период t находится в будущем. Может показаться странным, что эта модель не пытается экстраполировать временной тренд в будущее, но важно иметь в виду, что базовые оценки используются только как компонент в более продвинутых моделях и их цель — устранить наблюдаемые тенденции и уточнить сигнал, а экстраполяцию можно выполнить по оставшимся частям модели.

Такой подход может давать хорошие результаты для элементов, но быть не эффективным для пользователей, по крайней мере по двум причинам. Во-первых, один пользователь, как правило, имеет гораздо меньше рейтингов, чем один элемент, поэтому смещение нельзя надежно оценить, даже для относительно больших временных интервалов. Во-вторых, предвзятость пользователей может меняться гораздо быстрее, чем популярность элементов, и потому требуется использовать еще меньшие временные интервалы. Эту задачу можно решить путем моделирования дрейфа смещения пользователя как простой функции вместо точечных оценок. Например, для моделирования дрейфа можно использовать следующую функциональную форму:

$$d_u(t) = \text{sgn}(t - t_u) |t - t_u|^\beta, \quad (5.147)$$

где $\text{sgn}(x)$ — функция знака, равная единице, когда x положительно, и -1 , когда x отрицательно, t_u — средняя дата рейтинга, и β — постоянный параметр, выбранный

с помощью перекрестной проверки. Грубо говоря, это линейная функция, которую можно изогнуть, задав параметр $\beta < 1$, как показано на рис. 5.25. Смещение пользователя с течением времени можно определить как

$$b_u(t) = b_u + w_u \cdot d_u(t), \quad (5.148)$$

где стационарная часть b_u и коэффициент масштабирования w_u определяются по данным для каждого пользователя u . Эту модель можно улучшить включением дополнительных членов, отражающих дневную изменчивость смещения пользователя, такую как количество рейтингов, присвоенных пользователем в определенный день.

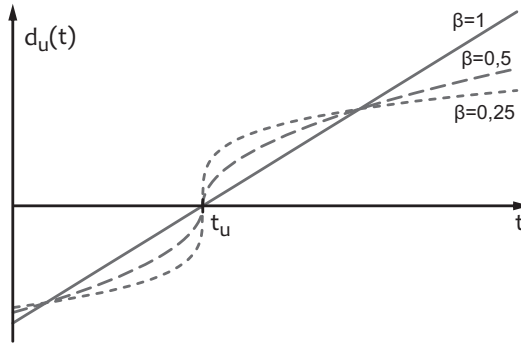


Рис. 5.25. Функция дрейфа смещения пользователя для разных значений β .
Функция линейна при $\beta = 1$

5.10.3.2. Модель ближайших соседей с затуханием во времени

Модель ближайших соседей прогнозирует рейтинг для заданной пары пользователь/элемент, усредняя рейтинги похожих пользователей или элементов. Рейтинги могут усредняться с использованием эвристической меры сходства или весовых коэффициентов, вычисленных по данным. Например, вот как можно определить одну из простейших моделей ближайших соседей на основе элементов (более полные и практичные версии этой модели рассматривались в разделе 5.7.4):

$$\hat{r}_{ui} = b_{ui} + \sum_{j \in Q_{ui}^k} w_{ij} (r_{uj} - b_{uj}), \quad (5.149)$$

где Q_{ui}^k — ближайшие соседи элемента i из множества элементов, оцененных пользователем, b_{uj} — базовая оценка, и вес w_{ij} — мера сходства элементов i и j , получен-

ная методом градиентного спуска, или некоторая эвристическая мера, такая как коэффициент корреляции Пирсона. Независимо от подхода к оценке веса можно утверждать, что, если пользователь присваивает одинаковые рейтинги двум элементам, i и j , это положительно влияет на соответствующий вес w_{ij} и в свою очередь повышает соответствующий рейтинг r_{ij} . Можно утверждать, что со временем рейтинги утрачивают релевантность, поэтому тот факт, что пользователю в прошлом одинаково нравились два элемента, не обязательно говорит о их сходстве в будущем из-за изменчивого характера вкуса и идентичности пользователя [Ding and Li, 2005; Koren, 2009]. Это обстоятельство можно учесть, добавив в модель коэффициент затухания во времени для девальвации старых оценок:

$$\hat{r}_{ui}(t) = b_{ui} + \sum_{j \in Q_{ui}^k} e^{-c_u |t - t_{uj}|} \cdot w_{ij} (r_{uj} - b_{uj}). \quad (5.150)$$

Скорость затухания c_u , как правило, зависит от пользователя и должна вычисляться как дополнительная переменная в процессе градиентного спуска.

5.10.3.3. Модель скрытых факторов с учетом временной динамики

Модели скрытых факторов предсказывают рейтинги, вычисляя корреляцию между представлениями пользователей и элементов в пространстве скрытых факторов. Эти модели можно расширить для учета временных эффектов, применив методы, разработанные выше, для базовых оценок. Напомню, что самая простая модель скрытых факторов определяется как

$$\hat{r}_{ui} = \mathbf{p}_u \mathbf{q}_i^T, \quad (5.151)$$

где \mathbf{p}_u и \mathbf{q}_i являются k -мерными векторами скрытых факторов пользователя и элемента соответственно. Изменчивость вкусов пользователей можно учесть, прибавив член временного дрейфа к каждому элементу вектора скрытых факторов, по аналогии с соответствующим выражением 5.148 для базовых оценок:

$$p_{us}(t) = p_{us} + w_{us} \cdot d_u(t), \quad 1 \leq s \leq k, \quad (5.152)$$

где p_{us} — стационарный член, $d_u(t)$ — дрейф, определяемый функцией 5.147, а s — индекс скрытого измерения. Обе части этой модели, стационарная и зависящая от времени, оцениваются по данным с использованием градиентного спуска.

На практике, однако, чаще используются более продвинутые модели скрытых факторов, чем базовое решение 5.151. К счастью, фактор, зависящий от времени

и определяемый выражением 5.152, можно использовать в большинстве моделей скрытых факторов. Особенно важной из них является модель SVD++, определяемая выражением 5.118. Вставив в модель SVD++ фактор, зависящий от времени, получаем модель под названием TimeSVD++:

$$\hat{r}_{ui} = \mu + b_i + b_u + \left(\mathbf{p}_u(t) + |\mathbf{I}_u|^{-\frac{1}{2}} \sum_{j \in \mathbf{I}_u} \mathbf{y}_j \right) \mathbf{q}_i^T. \quad (5.153)$$

Модель TimeSVD++ является одной из наиболее точных негибридных моделей совместной фильтрации и часто считается вершиной инженерии рекомендательных систем [Koren, 2009].

5.11. Неперсонализированные рекомендации

Качество рекомендаций, как правило, определяется способностью рекомендательной системы распознать намерение пользователя и найти предложения, соответствующие этому намерению. Вследствие широкого разнообразия намерений пользователей рекомендательным системам может сильно помочь возможность доступа и использования личной и поведенческой информации; системы, не имеющие доступа к таким данным и не использующие их, почти наверняка будут давать рекомендации низкого качества. Проще говоря, рекомендации, созданные для всех, то есть для абстрактного среднестатистического пользователя, едва ли пригодятся большинству, потому что лишь несколько реальных пользователей будут точно соответствовать среднему профилю. Однако рекомендательные системы должны учитывать, что доступная личная и поведенческая информация может быть неполной или ненадежной. Например, онлайн-системы часто имеют дело с анонимными пользователями, не имеющими или имеющими очень ограниченную историю взаимодействий, собранную во время текущего сеанса. Такие случаи очень распространены, поэтому важно расширить диапазон рекомендательных алгоритмов с помощью неперсонализированных методов, которые можно использовать в автономном режиме или объединять в гибридную систему с персонализированными методами. Точность неперсонализированных рекомендаций часто ниже, чем персонализированных, но они все еще могут быть эффективным решением в некоторых ситуациях, таких как перекрестные продажи товаров.

5.11.1. Типы неперсонализированных рекомендаций

Неперсонализированные методы рекомендаций обычно можно рассматривать как крайний случай контекстных рекомендаций, когда сигналы о целевом пользователе

полностью отсутствуют и рекомендуемые элементы выбираются исключительно на основе контекстной и справочной информации. Контекст, однако, может включать сведения о поведении других пользователей, поэтому общие шаблоны взаимодействия пользователь/элемент и статистика популярности элемента могут быть известны модели. Рассмотрим несколько примеров неперсонализированных рекомендаций, часто используемых на практике.

ПОПУЛЯРНЫЕ ЭЛЕМЕНТЫ. Популярные категории и бренды часто выделяются в пользовательских интерфейсах для упрощения навигации, а популярные товары часто продвигаются в таких разделах, как *Лидеры продаж* или *Топ-10*. Рекомендации этого типа могут не использовать контекст уровня запроса, но опираются на динамически обновляемую статистику продаж или просмотров, которую можно считать базовым контекстом. С точки зрения предиктивного моделирования эти методы просто используют тот факт, что наиболее частые решения о покупке являются лучшими прогнозами для целевого пользователя, когда дополнительная информация неизвестна.

НАБИРАЮЩИЕ ПОПУЛЯРНОСТЬ ЭЛЕМЕНТЫ. Рекомендательная система может рекомендовать продукты, демонстрирующие тенденцию к увеличению популярности, а не лидеров продаж, исходя из предположения, что такие рекомендации чаще оказываются приятными неожиданностями. Этот подход, в частности, способен лучше продвигать товары с длинным хвостом, потому что даже для медленно продвигаемого продукта могут случаться всплески популярности благодаря рекламе или активности в социальных сетях. Модель рекомендации элементов, набирающих популярность, обычно оценивает их, опираясь на сглаженную версию истории изменения объема продаж. Например, функцию оценки для элемента i можно определить как

$$s(i) = 1,00 \cdot \Delta v_1(i) + 0,50 \cdot \Delta v_2(i). \quad (5.154)$$

где $\Delta v_1(i)$ — относительное изменение объема продаж (в процентах) за предыдущий день, а $\Delta v_2(i)$ — изменение объема два дня назад.

НОВЫЕ ПОСТУПЛЕНИЯ. Некоторые рекомендательные системы выделяют и продвигают новые элементы или элементы, недавно добавленные в ассортимент.

ПОХОЖИЕ ЭЛЕМЕНТЫ. В маркетинговых онлайн-каналах неперсонализированные рекомендации часто могут создаваться на основе контекста просмотров. Типичным примером могут служить такие рекомендации, как *Похожие товары* или *Вам могут понравиться*, отображаемые на страницах с описаниями продуктов, то есть в контексте определенного продукта. Такие рекомендации часто можно генерировать с помощью стандартных методов поиска. Например, можно опре-

делить расстояние между двумя продуктами как взвешенное среднее расстояний $TF \times IDF$ между соответствующими полями документов продуктов и рекомендовать наиболее похожие товары. Эта стратегия основана на содержимом элемента и, следовательно, имеет тенденцию давать тривиальные и ожидаемые рекомендации.

ЧАСТО ВСТРЕЧАЮЩИЕСЯ ЗАКОНОМЕРНОСТИ. Рекомендации лидеров продаж используют только небольшую часть информации, доступной в истории покупок и просмотров, а именно общую статистику продаж. Контекстная информация, такая как продукт или категория, просматриваемые в данный момент, также не используется, что отрицательно сказывается на релевантности рекомендаций для перекрестных продаж. Более целенаправленные рекомендации можно генерировать путем анализа закономерностей покупок и просмотров и выявления товаров, которые часто покупаются вместе с данным товаром или набором товаров. Рекомендации этого типа часто присутствуют на страницах сведений о продуктах в таких разделах, как *Клиенты, купившие этот товар, также покупают* или *Товары, которые клиенты покупают после просмотра этого товара*. Мы подробно обсудим этот подход в следующем разделе.

Несмотря на то что вышеперечисленные методы классифицируются как неперсонализированные, важно понимать, что большинство из них способно достичь некоторого уровня персонализации путем сегментации и согласования с отдельными аспектами контекста. Например, рекомендательная система новостей может предложить пользователю выбрать темы, которые его интересуют (например, политика, наука, спорт и т. д.), а затем рекомендовать наиболее популярные или набирающие популярность элементы в указанных категориях вместо того, чтобы генерировать рекомендации на основе глобальной статистики. Рекомендательная система также может фильтровать или повторно ранжировать персонализированные рекомендации на основе статистики популярности или дат выпуска.

5.11.2. Рекомендации с использованием ассоциативных правил

Система неперсонализированных рекомендаций может анализировать сделки, совершенные в прошлом по различным маркетинговым каналам, включая обычные и онлайн-магазины, для выявления типичных зависимостей между продуктами, которые можно использовать для выработки рекомендаций. Например, если два элемента часто покупаются вместе, это может означать, что второй элемент является разумной рекомендацией с целью продвижения перекрестных продаж для пользователей, просматривающих первый элемент, и наоборот. Рекомендации создаются в контексте отдельного элемента (например, когда пользователь просматривает страницу с описанием определенного продукта) или нескольких элементов

(например, когда пользователь уже добавил несколько продуктов в корзину). Важно отметить, что граница между персонализированными и неперсонализированными рекомендациями в этом случае довольно размыта. Например, рекомендации, сгенерированные в контексте определенного продукта и показанные на странице с его описанием, можно рассматривать как неотъемлемую статическую часть этой страницы. Они одинаковы для всех пользователей и, следовательно, неперсонализированы. Однако если один и тот же контекст с единственным элементом постоянно привязан к пользователю и описывает его поведение, такой контекст можно рассматривать как персонализацию. Контекст, включающий несколько элементов, определенно можно интерпретировать как историю взаимодействий или неявную обратную связь. В этом случае рекомендации, сгенерированные на основе закономерностей, выявленных в данных о сделках, определенно можно классифицировать как совместную фильтрацию.

Если целью является получение рекомендаций на основе текущего просматриваемого элемента или нескольких элементов, для нас особенно важно будет выявить закономерности в исторических данных в форме следующих правил:

если пользователь приобретает элементы $X = \{i_1, i_2, \dots\}$,
тогда он приобретет также элементы $Y = \{j_1, j_2, \dots\}$.

Множества элементов X и Y называются *наборами элементов*, а ассоциативное правило, описанное выше, обозначается как $X \rightarrow Y$. Например, правило $\{\text{макароны, вино}\} \rightarrow \{\text{чеснок}\}$ указывает, что люди, покупающие макароны и вино вместе, могут также купить чеснок. Количество ассоциативных правил, хоть как-то подтверждаемых данными, то есть когда имеется хотя бы один случай одновременной покупки наборов элементов X и Y , может быть очень большим. Однако цель системы рекомендаций — найти правила, соответствующие устойчивым закономерностям, которые можно использовать в качестве признаков, предсказывающих поведение пользователей. Для выбора таких правил требуется ввести более формальные показатели качества.

Допустим, у нас есть массив сделок T , в котором каждая сделка представлена коллекцией элементов, приобретенных вместе. *Поддержку* (support) набора элементов X можно определить как долю сделок, содержащих все элементы из набора, то есть эмпирическую вероятность X :

$$\text{support}(X) = \frac{|t : X \subseteq t|}{|T|}, \quad t \in T. \quad (5.155)$$

Поддержка ассоциативного правила — это доля сделок, содержащих оба набора элементов из правила:

$$\text{support}(X \rightarrow Y) = \text{support}(X \cup Y). \quad (5.156)$$

Высокий уровень поддержки гарантирует, что правило соответствует устойчивой закономерности совместной покупки наборов элементов. Однако может случиться так, что эти наборы элементов часто покупаются отдельно, то есть высокая поддержка лишь подтверждает популярность обоих наборов, тогда как в действительности между ними нет никакой зависимости. Этот аспект измеряется *достоверностью* (confidence) правила, определяемой как доля сделок, содержащих X , которые содержат также Y :

$$\text{confidence}(X \rightarrow Y) = \frac{\text{support}(X \cup Y)}{\text{support}(X)}. \quad (5.157)$$

Достоверность можно интерпретировать как условную вероятность обнаружить набор элементов Y в сделке, содержащей также набор X , то есть $\Pr(Y|X)$. Обратите внимание, что поддержка и достоверность определяются на основе вероятности покупки и, соответственно, могут быть связаны с денежными мерами, такими как доход [Ju et al., 2015; Geng and Hamilton, 2006]. Например, ожидаемую доходность (revenue) правила (и рекомендации, сгенерированной на основе этого правила) можно приблизительно оценить следующим образом:

$$\text{revenue}(X \rightarrow Y) = \text{support}(X \rightarrow Y) \cdot \sum_{i \in Y} \text{price}(i), \quad (5.158)$$

если учесть, что пользователь собирается купить элемент X . Более точные оценки денежных метрик можно получить с помощью методов моделирования подъема, рассмотренные выше в контексте оптимизации продвижения. Рекомендательной системе обычно нужны ассоциативные правила с высоким уровнем поддержки и достоверности, гарантирующим их надежность и избирательность. Создание таких правил на основе имеющейся истории сделок является стандартной задачей анализа данных, известной как частотный анализ набора элементов, анализ сходства или анализ покупательского поведения (анализ потребительской корзины). Эту задачу можно решить с помощью широкого спектра специализированных алгоритмов, таких как Apriori или FP-growth.

ПРИМЕР 5.7

Рассмотрим пример, иллюстрирующий использование ассоциативных правил для создания рекомендаций. В отличие от традиционной совместной фильтрации, для выявления ассоциативных правил требуют более детализированные данные уровня сделок, но сделки не должны быть свя-

заны с конкретными пользователями (то есть неважно, какой пользователь выполнил ту или иную сделку). Мы проанализируем выборку из истории сделок продуктового магазина, представленную в табл. 5.9.

Таблица 5.9. Пример истории сделок

№ сделки	Элементы
1	молоко, хлеб, яйца
2	хлеб, сахар
3	молоко, хлопья
4	хлеб, хлопья
5	молоко, хлеб, сахар
6	хлопья, молоко, хлеб
7	хлеб, хлопья
8	молоко, хлопья
9	молоко, хлеб, хлопья, яйца

Чтобы сгенерировать рекомендации, например, для молока, применим алгоритм поиска ассоциативных правил к истории транзакций с ограничением $X = \{\text{молоко}\}$ и отсортируем правила по степени достоверности. Результат показан в табл. 5.10. Степень достоверности для правила *молоко* \rightarrow *хлопья*, например, равен $4/6$, потому что четыре сделки содержат оба этих элемента и шесть сделок содержат молоко. Рекомендуемые элементы извлекаются из правых частей правил, поэтому в список рекомендаций для молока попадают хлопья, хлеб, яйца и сахар, в порядке релевантности.

Таблица 5.10. Ассоциативные правила для молока

Ранг	Правило	Поддержка	Достоверность
1	молоко \rightarrow хлопья	$4/9$	$4/6$
2	молоко \rightarrow хлеб	$4/9$	$4/6$
3	молоко \rightarrow яйца	$2/9$	$2/6$
4	молоко \rightarrow сахар	$1/9$	$1/6$

Обратите внимание: даже притом что анализ покупательского поведения является методом обучения без учителя, в этой конкретной ситуации мы

фактически решаем задачу классификации и выбора признаков, потому что контекст (молоко в данном примере) можно рассматривать как обучающую метку, а другие элементы в сделках — как признаки, и цель состоит в том, чтобы определить признаки с наибольшей прогнозирующей способностью. В подходе на основе ассоциативных правил эти признаки соответствуют правой части правил с высоким уровнем достоверности.

Как отмечалось в начале этого раздела, ассоциативные правила также можно использовать для персонализированных рекомендаций, в зависимости от количества элементов, включенных в контекст (левая часть правила) и семантического значения контекста. Для таких случаев использования, как персонализация анонимного веб-сеанса, применение ассоциативных правил может оказаться эффективной альтернативой другим методам совместной фильтрации, таким как модели ближайших соседей, как в смысле точности, так и вычислительной сложности [Mobasher et al., 2001].

5.12. Многоцелевая оптимизация

Все методы рекомендаций, описанные выше, в основном преследуют одну цель — обеспечить наилучшее семантическое соответствие или прогнозируемую оценку предпочтений. Однако точность рекомендаций может быть не единственной целью при разработке рекомендательной системы: маркетолог также может захотеть включить нескольких конкурирующих целей в рекомендации, предлагаемые клиентам. Например, продуктовые магазины могут быть заинтересованы в продвижении скоропортящихся продуктов с более коротким сроком хранения, магазины модной одежды — в продвижении брендов спонсоров или сезонных коллекций, а розничные торговцы в целом могут извлечь выгоду из рекомендаций продуктов с более высокой доходностью или, учитывая ограниченность складских помещений, захотеть избежать затоваривания складов [Jambor and Wang, 2010].

Один из возможных подходов к реализации многоцелевой рекомендательной системы заключается в смешивании сигналов семантической релевантности с сигналами, соответствующими вторичным целям. С этой точки зрения многоцелевые методы можно сравнить с гибридными моделями, смешивающими несколько сигналов для достижения оптимальных результатов. Основное отличие заключается в том, что в качестве цели оптимизации гибридные методы обычно используют стандартные функции потерь, такие как средняя ошибка прогнозирования рейтинга, тогда как многоцелевые модели используют более специализированные цели оптимизации. В этом разделе мы рассмотрим рекомендательную систему

с несколькими целями, которая была разработана и проверена на практике социальной сетью LinkedIn, ориентированной на поиск работы [Rodriguez et al., 2012]. Для LinkedIn основной целью является рекомендация кандидатов, которые семантически соответствуют описанию работы, а вторичной целью — демонстрация поведения, связанного с поиском работы.

Начнем с идеи повторного ранжирования рекомендаций, полученных с помощью основного алгоритма, с применением функции, оптимизирующей некоторую вторичную цель, с условием штрафования отклонений от исходного ранжирования на основе релевантности. Рассмотрим сначала случай с одной вторичной целью и определим относительно абстрактную основу, которую можно адаптировать для широкого круга целей. Во-первых, основной алгоритм рекомендаций можно использовать для ранжирования всех m элементов для каждого из n пользователей. Обозначим эти исходные рекомендации как $n \times m$ матрицу \mathbf{Y} , в которой строки соответствуют пользователям, столбцы — элементам, а каждое значение в матрице представляет ранг элемента в списках рекомендаций. Предположим, что каждый пользователь получает только $k \ll m$ рекомендаций, но при этом оценке подвергаются все элементы, чтобы функция повторного ранжирования имела достаточно широкий выбор. На практике необязательно ранжировать все m элементов — количество рекомендаций можно ограничить некоторым числом, которое существенно больше k . Затем каждый рекомендуемый элемент можно оценить в соответствии со вторичной целью. Обозначаем матрицу $n \times m$ этих оценок как \mathbf{X} . Эта матрица, например, может содержать валовую прибыль продукта. Обратите внимание, что оценка может быть функцией как самого элемента, так и его положения в матрице рекомендаций \mathbf{Y} , то есть пользователя и ранга. Задачу оптимизации тогда можно определить следующим образом:

$$\begin{aligned} \max_w \quad & g(\phi(\mathbf{Y}, \mathbf{X}, w)) \\ \text{при условии} \quad & d(\text{top}_k(\mathbf{Y}), \text{top}_k(\phi(\mathbf{Y}, \mathbf{X}, w))) \leq c, \end{aligned} \quad (5.159)$$

где

- g — функция полезности, которая оценивает качество рекомендации с точки зрения вторичной цели;
- ϕ — составная функция ранжирования, объединяющая пары строк из матриц \mathbf{X} и \mathbf{Y} в новый список рекомендаций, уравнивающий две цели;
- w — параметр (вес), определяющий баланс смешивания двух целей. Данный параметр является предметом оптимизации;
- $\text{top}_k(\cdot)$ — обозначает первые k элементов с максимальными значениями рангов. Эта операция усекает исходные матрицы \mathbf{X} и \mathbf{Y} до размера $n \times k$;

- d — функция расстояния, измеряющая расхождение между двумя матрицами рекомендаций. Один из возможных способов измерить расхождение между двумя векторами оценок \mathbf{x} и \mathbf{y} — вычислить сумму квадратов ошибок между их гистограммами:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^b \left(H(\mathbf{x})_i - H(\mathbf{y})_i \right)^2, \quad (5.160)$$

где гистограмма $H(\mathbf{x})$ — это вектор с b элементами (сегментами), каждому элементу которого соответствует количество оценок в \mathbf{x} , попадающих в соответствующий диапазон. Общее расстояние между матрицами определяется как сумма расстояний между всеми пользователями, то есть строками матриц;

- c — пороговое значение, ограничивающее величину расхождения между первоначальным и ранжированным списками рекомендаций.

Основная идея оптимизации заключается в увеличении полезности переупорядоченных рекомендаций смешиванием оценок релевантности со второстепенной целью, но с наложением штрафа за разницу между первоначальным порядком следования релевантных рекомендаций и порядком после повторного ранжирования, чтобы гарантировать, что релевантность не будет принесена в жертву вторичной цели. Функция ϕ должна включать настраиваемые параметры, которые управляют балансом между двумя целями и являются предметом оптимизации. Этот подход легко можно распространить на случай с большим количеством целей и ограничений расхождения. Обозначив число целей как q , можно определить следующую задачу многоцелевой оптимизации:

$$\begin{aligned} \max_{\mathbf{w}} \quad & g(\phi(\mathbf{Y}, \mathbf{X}, \mathbf{w})) \\ \text{при условии} \quad & d_j(\text{top}_k(\mathbf{Y}), \text{top}_k(\phi(\mathbf{Y}, \mathbf{X}, \mathbf{w}))) \leq c_j, \end{aligned} \quad (5.161)$$

где \mathbf{X} теперь $n \times m \times q$ матрица оценок, \mathbf{w} — вектор q весовых параметров, а j перебирает все критерии расхождений.

Проиллюстрируем, как представленную выше модель оптимизации можно адаптировать к практическим задачам, на нескольких примерах. Сначала рассмотрим ретейлера, который хочет включить в оценки рекомендаций цель увеличения доходности. Общую функцию полезности можно определить как ожидаемую валовую прибыль в предположении, что $M(i) \in [0, 1]$ является нормализованной валовой прибылью для элемента i , а вероятность покупки моделируется как величина, обратная позиции в рейтинге (то есть чем ниже элемент в списке рекомендаций, тем ниже вероятность его покупки):

$$g(\mathbf{z}) = \frac{1}{k} \sum_{i=1}^m \frac{M(i)}{z_i} \cdot \mathbb{I}(z_i \leq k), \quad (5.162)$$

где \mathbf{z} — вектор рангов, вычисляемых составной функцией ранжирования ϕ , а \mathbb{I} — индикаторная функция, равная единице для истинного аргумента, и нулю в противном случае. Поскольку вторичной целью является ожидаемая валовая прибыль, матрица \mathbf{X} прямо определяется как

$$x_{ii} = M(i). \quad (5.163)$$

Составную функцию ранжирования в этом случае можно определить как сочетание исходной оценки релевантности y , возвращаемой основным алгоритмом рекомендации, и оценки доходности x :

$$\mathbf{z} = \phi(\mathbf{y}, \mathbf{x}): z_i = y_i \cdot x_i^w, \quad (5.164)$$

где w — параметр, управляющий балансом между релевантностью и подъемом продуктов с высокой доходностью. Этот параметр является предметом оптимизации в задаче 5.159.

Второй пример повторного ранжирования в соответствии со вторичной целью — это подъем рекламируемых элементов, таких как имеющиеся в продаже или скоропортящиеся продукты. Функцию полезности можно определить как среднее количество рекламируемых продуктов в окончательном списке k рекомендаций:

$$g(\mathbf{z}) = \frac{1}{k} \sum_{i=1}^m F(i) \cdot \mathbb{I}(z_i \leq k), \quad (5.165)$$

где $F(i)$ — метка, равная единице, если элемент является рекламируемым, и нулю в противном случае. Матрица \mathbf{X} определяется как

$$x_{ii} = F(i). \quad (5.166)$$

Составная функция ранжирования объединяет оценку релевантности и метку с параметром балансировки w , который является предметом оптимизации:

$$\mathbf{z} = \phi(\mathbf{y}, \mathbf{x}): z_i = y_i \cdot w^{x_i}. \quad (5.167)$$

Эту функцию ранжирования легко можно распространить на случай с несколькими метками, каждая из которых вносит свой вклад в окончательный ранг, в соответствии со своим собственным параметром баланса (напомню, что \mathbf{X} — это матрица оценок $n \times m \times q$, соответственно \mathbf{x} — матрица $q \times m$):

$$\mathbf{z} = \phi(\mathbf{y}, \mathbf{x}): z_i = y_i \cdot w_1^{x_{1i}} \cdot w_2^{x_{2i}} \cdot \dots \cdot w_q^{x_{qi}}. \quad (5.168)$$

Задача оптимизации 5.159 зависит от ранжирующей функции, поэтому стандартные методы оптимизации для гладких функций, такие как градиентный спуск,

не применимы напрямую. В общем случае эту задачу можно решить с помощью алгоритмов обучения ранжированию [Rodriguez et al., 2012]. Однако во многих практических приложениях используется только один или два параметра w . В этом случае задачу можно решить простым перебором всех возможных значений.

5.13. Архитектура систем рекомендаций

До сих пор мы рассматривали широкий спектр моделей и алгоритмов рекомендаций, а также методы объединения нескольких моделей в гибридные или корректировки рекомендаций на основе контекстной информации или вторичных целей. Однако система рекомендаций — это нечто большее, чем просто реализация некоторого алгоритма. Это сложная программная система, включающая несколько компонентов и модулей, которые связывают модель рекомендаций с внешним миром и обеспечивают ее функционирование. В этом разделе мы обсудим возможную эталонную архитектуру системы рекомендаций, представленную на рис. 5.26, а также некоторые возможные варианты и компромиссы [Jack et al., 2016].

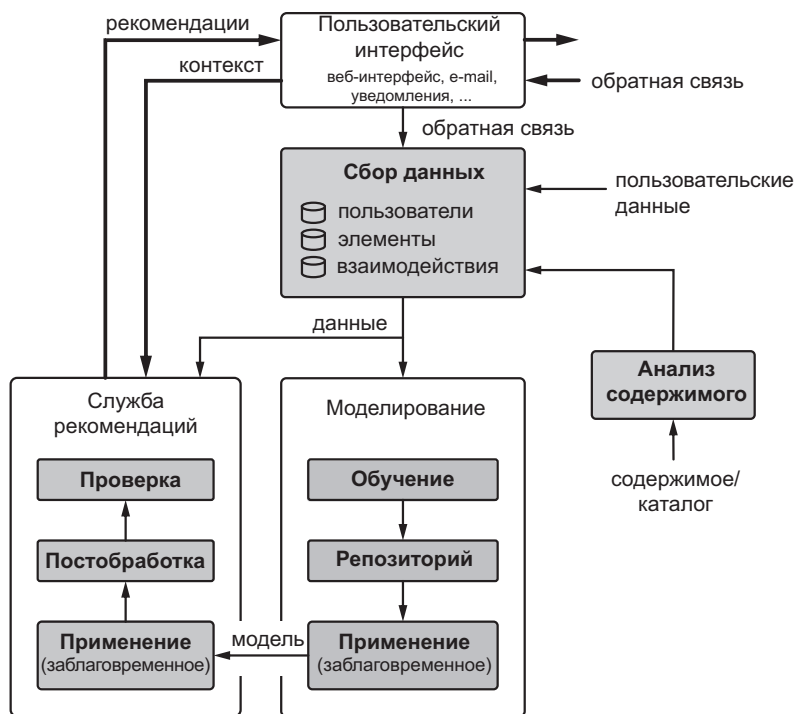


Рис. 5.26. Обобщенная архитектура рекомендательной системы

ПОЛЬЗОВАТЕЛЬСКИЙ ИНТЕРФЕЙС. К рекомендательной системе может быть подключено несколько пользовательских интерфейсов, таких как веб-сайты, электронная почта, мобильные уведомления или ленты новостей. Эти каналы взаимодействуют с основной рекомендательной системой через интерфейс, получающий контекстную информацию в виде аргумента и возвращающий ранжированные рекомендации. В большинстве случаев интерфейс рекомендательной системы может быть таким же простым, как боковая панель на веб-странице, содержащая рекомендации. Службы, сильно зависящие от рекомендаций, например онлайн-видеосервисы, часто предоставляют гораздо более полные интерфейсы, включающие несколько разделов, например с персонализированными рекомендациями, популярными элементами и последними тенденциями.

СБОР ДАННЫХ. Общее число источников данных, используемых рекомендательной системой, может быть очень большим. Одна из причин этого заключается в том, что промышленные рекомендательные системы часто используют множество разных алгоритмов либо для создания гибридных моделей, либо для экспериментов, соответственно система должна иметь в своем распоряжении достаточный объем данных для поддержки фильтрации по содержанию, совместной фильтрации и методов, основанных на популярности. Кроме того, система может иметь доступ к некоторым внешним или сторонним данным или сигналам, помогающим уточнить рекомендации. Это часто требует создания обширной инфраструктуры сбора данных, объединяющей профили пользователей (например, предпочтения и личную информацию); сведения об элементах, как правило, требующих применения некоторого механизма, выполняющего анализ и преобразование исходных данных; данные о взаимодействиях пользователь/элемент, которые обычно поступают через пользовательский интерфейс в виде рейтингов и других видов обратной связи.

МОДЕЛИРОВАНИЕ. Собранные, отфильтрованные и объединенные данные используются для создания моделей рекомендаций. Модели развертываются в репозитории и периодически проходят повторное обучение для учета изменений в данных. Модель может применяться в процессе моделирования (заблаговременно), или эксплуатации (оперативно), или на обоих этапах. Если модель применяется исключительно заблаговременно, результатом является набор рекомендаций для всех пользователей. Рекомендации могут обновляться и загружаться в службу рекомендаций по расписанию, например ежедневно. Если модель частично применяется в оперативном режиме, заблаговременно вычисляются только определенные элементы данных, такие как матрицы сходств элементов или векторы скрытых факторов. Такой подход позволяет выполнить тяжелые вычисления заблаговременно и сохранить гибкость оперативных рекомендаций. Промышленные рекомендательные системы обычно имеют репозиторий с несколькими алгоритмами рекомендаций, которые постоянно обновляются и тестируются.

СЛУЖБА РЕКОМЕНДАЦИЙ. Основная цель службы рекомендаций — применение модели и передача рекомендаций пользовательским интерфейсам. Служба рекомендаций может реализовать ряд функций по согласованию с контекстом и оперативному управлению. Во-первых, рекомендации, полученные основным алгоритмом, могут подвергаться постобработке для внесения дополнительных улучшений и корректировок. Например, служба может в режиме реального времени отслеживать рекомендации, которые пользователь уже видел, и прокручивать или произвольно перемешивать список рекомендаций, чтобы создать у пользователя ощущение динамики, продуктивности и приятной неожиданности. Во-вторых, служба может контролировать качество рекомендаций с помощью правил проверки и журнала ошибок или автоматически вносить корректировки, если проверка не увенчалась успехом. Примерами автоматического контроля качества могут служить проверка общего количества рекомендаций в списке и мониторинг времени выполнения алгоритма.

5.14. Итоги

- Цифровые каналы позволяют маркетологам продвигать широкие и глубокие ассортименты с большим количеством нишевых продуктов. Это одно из основных их отличий от традиционных каналов, где ассортимент ограничен издержками распространения.
- Чрезвычайно широкий и глубокий ассортимент с длинным хвостом нишевых продуктов порождает потребность в эффективных услугах, помогающих пользователям обнаружить их, включая поиск и рекомендации.
- Основная цель рекомендательных систем — дать клиентам релевантные предложения в случаях, когда намерение совершить покупку не выражено явно. Их можно рассматривать как дополнение к службам поиска, где намерение явно выражено в поисковом запросе.
- Рекомендательные системы обычно могут использовать данные о взаимодействиях пользователей с элементами, включая явно присвоенные рейтинги и неявно собранные истории просмотров, данные каталогов и контекстную информацию. Основным результатом работы службы рекомендаций является ранжированный список рекомендуемых элементов.
- Одними из основных входных данных рекомендательной системы являются рейтинги клиентов. Рейтинги обычно представлены в виде матрицы, строки которой соответствуют клиентам, столбцы — элементам, а числовые значения — присвоенным рейтингам. Значения рейтингов также могут сопровождаться контекстной информацией, такой как метки времени или маркетин-

говые каналы, из которых были получены рейтинги. Матрицы рейтингов обычно очень разрежены.

- Основными бизнес-целями рекомендательных систем являются релевантность, новизна, неожиданность и разнообразие рекомендаций. Релевантность рекомендаций можно измерить с помощью точности прогнозирования рейтинга и точности ранжирования. Также можно определить количественные показатели для новизны, неожиданности и разнообразия.
- Наиболее важными семействами алгоритмов рекомендаций являются фильтрация по содержанию и совместная фильтрация. Эти основные алгоритмы можно расширить путем гибридизации, согласования с контекстом и добавления дополнительных целей и сигналов.
- Фильтрация по содержанию рекомендует элементы, похожие на элементы, понравившиеся пользователю в прошлом. Фильтрацию по содержанию можно рассматривать как задачу классификации элементов. К ключевым преимуществам фильтрации по содержанию относятся: возможность рекомендовать элементы только на основе собственных рейтингов пользователя и новые элементы без рейтинга, а также интерпретируемость результатов. В числе основных недостатков можно назвать комплексное проектирование признаков, необходимое для классификации контента, и смещенность в сторону тривиальных рекомендаций.
- В фильтрации по содержанию можно использовать меры сходства документов и другие методы поиска, такие как латентно-семантический анализ и латентное распределение Дирихле, для выбора наиболее похожих элементов. Альтернативное решение заключается в использовании моделей классификации текста, таких как наивный байесовский классификатор.
- Для поиска элементов или пользователей с похожими закономерностями в рейтингах совместная фильтрация использует матрицу рейтингов. Этот метод, как правило, генерирует более разнообразные и неожиданные рекомендации, чем фильтрация по содержанию. Совместную фильтрацию можно рассматривать как задачу восстановления матрицы.
- К наиболее важным методам совместной фильтрации относятся модели ближайших соседей и скрытых факторов.
- Несколько моделей рекомендаций можно объединить в одну гибридную модель. В гибридных моделях используются те же методы, что и в конвейерах смешивания сигналов в службах поиска. Гибридная модель может переключаться между сигналами релевантности, поступающих из составляющих ее моделей, смешивать их или передавать результаты одной модели на вход другой.

- Контекстные модели используют дополнительные атрибуты, такие как время или местоположение, чтобы сделать рекомендации более целенаправленными. Эти атрибуты можно рассматривать как дополнительные измерения, превращающие матрицу рейтингов в многомерный куб. Рекомендательная система может использовать контекстную информацию для предварительной фильтрации входных данных, заключительной фильтрации рекомендаций или расширения входных данных модели прогнозирования рейтингов. Контекстные модели используют те же идеи, что и метод контролируемого снижения точности в службах поиска.
- Рекомендации для неизвестных пользователей и пользователей с ограниченными историями взаимодействий и профилями являются особым случаем. Для получения неперсонализированных или частично персонализированных рекомендаций рекомендательные системы могут использовать базовую статистику продаж (лидеры продаж), содержимое (похожие товары) и закономерности покупок (часто выбираемые наборы элементов).
- К основным компонентам рекомендательных систем относятся: пользовательский интерфейс, слой сбора данных, слой моделирования и служба рекомендаций.

6

Ценообразование и ассортимент

Задача управления ценами имеет очень долгую историю. Фундаментальные аспекты ценообразования изучались веками, чтобы объяснить взаимовлияние спроса и предложения на рынке. Это привело к разработке комплексной теории, описывающей стратегические аспекты ценообразования, такие как структура цен, зависимость между ценой и спросом и другие. Эта теория предлагает относительно грубые методы оптимизации, которые, однако, способствуют формированию стратегических решений ценообразования. Возможность автоматического совершенствования тактических решений была впервые признана и использована в авиационной отрасли в начале 1980-х и отчасти объясняется развитием цифровых систем бронирования, которые обеспечили динамичное и гибкое управление ресурсами и ценами. Это потребовало разработки набора совершенно новых методов оптимизации, которые позже были заимствованы другими сферами услуг, такими как гостиницы и аренда автомобилей. Этот новый, по-настоящему алгоритмический подход, который обычно называют *управлением доходами* или *управлением доходностью*, наглядно продемонстрировал мощь автоматизированного управления ценами и ресурсами в многочисленных случаях, когда не успевшие внедрить новые методы обанкротились или проиграли конкурентную борьбу пионерам автоматизированного управления ценами.

Управление ценами тесно связано с другими программными услугами, особенно с продвижением товаров и рекламой. Методы управления ценами можно использовать как для оптимизации скидок во время рекламных акций, так и для формирования цен на рекламные и медийные ресурсы, продаваемые клиентам службы. Эту главу мы начнем с обзора основных принципов стратегического ценообразования и оптимизации цен. Затем займемся разработкой более тактических и практических методов прогнозирования спроса и оптимизации цен для

сегментации рынка, уценок и распродаж. Мы также кратко рассмотрим основные методы распределения ресурсов, используемые в сфере услуг для установления лимитов бронирования. Наконец, рассмотрим задачу оптимизации ассортимента, для решения которой можно повторно использовать некоторые строительные блоки, разработанные для управления ценами.

6.1. Среда

Как будет показано далее в этой главе, управление ценами является критически важным фактором, определяющим прибыльность предприятия и, в конечном счете, его выживание. Следовательно, процессы управления ценами в реальной жизни часто включают несколько уровней принятия решений, от стратегических административных решений до микрорешений на уровне отдельных сделок. Мы спрячем все эти сложности за фасадом относительно простой модели, представленной на рис. 6.1. Эта модель не учитывает явно некоторые стратегические аспекты ценообразования, но отражает основные особенности ценовой среды, важные для принятия микрорешений:

- Обычно предполагается, что компания продает некоторые продукты своим клиентам и на каждом продукте i зарабатывает прибыль G :

$$G_i = Q_i (P_i - V_i) - C_i, \quad (6.1)$$

где Q_i — проданное количество, P_i — цена, V_i — переменные стоимости (например, оптовая цена продукта), C_i — постоянные затраты, связанные с продуктом, и i — индексы продуктов в ассортименте, предлагаемых клиентам. Большинство методов, рассматриваемых в этой главе, ориентированы на максимизацию прибыли G как функции цены, хотя также важно иметь в виду, что эта оптимизация может быть предметом внешних стратегических ограничений. Например, компания может выбрать стратегию, ориентированную на конкуренцию и захват доли рынка, а не на прибыль, и ограничивать процесс оптимизации ценовой политикой.

- Прибыль является функцией количества проданного товара, которое, в свою очередь, зависит от спроса. Ключевым допущением в модели является неоднородность спроса, то есть изменчивость спроса по одному или нескольким измерениям, таким как сегменты клиентов, местоположение магазинов, сезоны, классы обслуживания и т. д. Это дает возможность дифференцировать цены по таким параметрам и, соответственно, регулировать прибыль на уровне от-

дельных клиентских сегментов или временных интервалов. Большую часть этой главы мы посвятим рассмотрению структуры и методов оптимизации цен с учетом различных аспектов неоднородности спроса.



Рис. 6.1. Среда управления доходностью

Спрос является функцией цены и других переменных, к которым можно отнести широкий спектр факторов, от цен у конкурентов до погоды. В самом простом случае система управления доходами может оптимизировать цены, создавая регрессионные модели потребности отдельных продуктов и определяя оптимальные цены, максимизирующие прибыль за счет увеличения спроса:

$$p_i^{opt} = \operatorname{argmax}_{p_i} Q_i(p_i) \cdot (p_i - V_i) - C_i, \quad (6.2)$$

где $Q_i(p_i)$ — модель прогнозирования спроса для продукта i . На практике эта задача, как правило, намного сложнее из-за различных ограничений и взаимозависимостей.

- Одним из примечательных примеров может служить ограничение уровня запасов — если продавец имеет ограниченный запас продукта, проданное количество Q является минимумом от спроса и уровня запасов.
- Еще одним важным фактором является зависимый спрос — поскольку продукты в одной категории часто взаимозаменяемы, изменение цены на один продукт может заставить клиентов переключиться на другой. Это усложняет задачу оптимизации, потому что цены на продукты в этом случае должны оптимизироваться совместно, а не по отдельности.
- Наконец, продавец может преследовать дополнительные цели, влекущие дополнительные ограничения. Например, розничный продавец модной одежды может стремиться продать товар к концу сезона.

Следовательно, система управления доходами имеет множество входных данных, включая историю изменения спроса, постоянные и переменные затраты, уровень запасов и другие бизнес-ограничения.

- Несмотря на то что оптимизация цен является естественной задачей для системы управления доходами, многие среды предлагают другие важные средства управления. Один из наборов таких средств управления связан с тем, как цена сообщается клиентам. Как мы обсудим позже, сообщать об изменениях цен зачастую эффективнее в форме скидок и специальных предложений, а не просто изменять базовые цены. Это устанавливает связь между задачами оптимизации цен и скидок. Второй набор средств управления связан с доступностью продукта и классами обслуживания. Классическим примером может служить авиакомпания, требующая бронировать недорогие билеты заранее, и делает эту возможность недоступной за несколько дней до вылета. Наконец, система управления доходами может управлять ассортиментом продукции, ее презентацией и вариантами размещения. Примерами могут служить оптимизация пространства на полках путем удаления плохо продаваемых продуктов и оптимизация внутренней планировки магазина путем размещения связанных продуктов вместе для увеличения перекрестных продаж.

Представленное описание модели определяет основные области для применения методов программного управления ценами. На следующих страницах мы познакомимся с основными понятиями и принципами ценообразования, а затем приступим к разработке методов оптимизации цен для среды, описанной в этом разделе.

6.2. Влияние ценообразования

Ценообразование играет критически важную роль в экономике предприятий, потому что цены являются ключевыми факторами, определяющими доходность и прибыльность. Правильные ценовые решения могут обеспечить значительное увеличение прибыльности, тогда как ошибочные могут иметь серьезные последствия. Одна из причин заключается в том, что цена определяет, как продукт или услуга позиционируется на рынке и воспринимается клиентами. Слишком низкие цены подрывают устойчивость фирмы, лишая ее недополученной прибыли и воспитывая у клиентов неправильные ожидания в отношении ценности и качества продукции; слишком высокие цены наносят вред продажам и репутации фирмы, что замедляет рост бизнеса. Другая причина заключается в сильной зависимости между ценами и прибылью в большинстве отраслей и предприятий, поэтому цена часто доминирует в уравнениях прибыли. Рассмотрим пример, иллюстрирующий важность ценовых решений.

ПРИМЕР 6.1

Рассмотрим воображаемого продавца одежды, который продает 100 000 предметов одежды ежемесячно по 40 долларов за штуку при оптовой цене 25 долларов за штуку и имеет постоянные затраты в размере 500 000 долларов в месяц. Прибыль этого продавца можно выразить в виде функции цены, затрат и объема продаж, используя базовое уравнение прибыли:

$$G = Q(P - V) - C, \quad (6.3)$$

где Q — проданное количество, P — цена, V — переменные затраты и C — постоянные затраты. То есть базовая прибыль будет:

$$G = 100,000 \times (\$40 - \$25) - \$500,000 = \$1,000,000. \quad (6.4)$$

Продавец может выбрать одну из нескольких стратегий для увеличения базовой прибыли. Одна из них — увеличение объема продаж путем инвестирования в маркетинговые кампании или в развитие новых каналов продаж. В числе других подходов можно назвать увеличение продажной цены, смену поставщика с целью уменьшения переменных затрат или сокращение постоянных затрат. Оценим все эти стратегии, как показано в табл. 6.1, и под-

считаем, как изменение объема продаж, цены, переменных и постоянных затрат на 1 % повлияет на прибыль. Как оказывается, прибыль наиболее чувствительна к изменению цены, что свидетельствует о высокой важности ценовых решений.

Таблица 6.1. Количественный пример, иллюстрирующий влияние цены, затрат и объема продаж на доходность

	Базовая стратегия	+1 % к объему продаж (Q)	+1 % к цене (P)	–1 % из переменных затрат (V)	–1 % из постоянных затрат (C)
Объем продаж (Q)	100 000	101 000	100 000	100 000	100 000
Цена (P)	40,00	40,00	40,40	40,00	40,00
Переменные затраты (V)	25,00	25,00	25,00	24,75	25,00
Постоянные затраты (C)	500 000	500 000	500 000	500 000	495 000
Прибыль (G)	1 000 000	1 015 000	1 040 000	1 025 000	1 005 000
ΔG %		+1,5 %	+4,0 %	+2,5 %	+0,5 %

В этом примере мы использовали довольно произвольные числа, однако эта модель преобладает в самых разных предприятиях во многих отраслях промышленности. Например, компания McKinsey and Associates провела исследования и, проанализировав прибыль 2 463 компаний, пришла к выводу, что изменение цены на 1 % дало в результате увеличение прибыли на 11,1 %, тогда как изменение объема продаж, переменных или постоянных затрат на 1 % дало увеличение на 3,3 %, 7,8 % и 2,3 % соответственно [Marn and Roseillo, 1992].

6.3. Цена и стоимость

Для разработки автоматизированных систем управления ценами мы должны разбить ценообразование на формальные задачи оптимизации, рассматривающие прибыль только как математические функции. С другой стороны, цена — сложный вопрос, зависящий от природы продукта, наличия конкуренции и психологии кли-

ента. В этом разделе мы начнем устранять разрыв между фундаментальной задачей ценообразования и задачами оптимизации, рассматривая цену и ее определение. Эта тема является стратегической и дает очень мало информации, как внедрить автоматизированное управление ценами, но содержит рекомендации, которые помогут разработать более сложные методы.

6.3.1. Ценовые границы

Как утверждает экономическая теория, цена определяется спросом и предложением на рынке. Каждый продукт или услуга имеет свою себестоимость, которую иногда можно рассматривать как «справедливую» базовую цену, однако ценообразование требует от нас углубиться в логику оценки ценности продавца и покупателя.

С одной стороны, можно предположить, что продавец товара или услуги имеет минимально выгодную цену. Продажа по цене выше этого значения приносит прибыль; продажа по цене ниже влечет убытки. Во многих случаях можно предположить, что эта базовая цена равна предельной себестоимости продукта.

С другой стороны, покупатель извлекает из приобретенного продукта определенную *пользу*. Полезность зависит от функциональных свойств продукта, способности клиента достичь полезных целей при использовании этих свойств, наличия продукта в нужном месте и в нужное время и других факторов. В некоторых случаях можно получить относительно точную оценку полезности. Например, полезность промышленного электрогенератора можно оценить по цене производимой им электроэнергии. В других случаях может достичь лишь грубого приближения к оценке полезности. Примером может служить инновационное медицинское оборудование, которое можно оценить столь же высоко, как человеческую жизнь.

Продавец и покупатель могут заключить взаимовыгодную сделку, если издержки окажутся ниже уровня полезности. В противном случае им обоим лучше отказаться от сделки. Следовательно, цена — это, по сути, диапазон, а не точка на числовой прямой. На очень высоком уровне целью системы управления ценами является динамический выбор оптимальной точки в этом диапазоне для отдельных сделок.

Издержки и полезность зачастую дают весьма широкие ценовые рамки, которые на практике могут оказаться бесполезными. Например, покупка бутылки содовой в жаркий день может быть буквально вопросом жизни и смерти, что резко повышает полезность, но высокая конкуренция удерживает цену ближе к нижней границе. С другой стороны, незначительные издержки на распространение программного обеспечения не препятствует удержанию цены вблизи верхней границы, определяемой полезностью. Более узкие рамки часто можно получить, сравнивая продукт или услугу с имеющимися альтернативами и тщательно оценивая их характеристики.

При наличии сопоставимой альтернативы ее цену можно принять за базовую. Цена на данный продукт может быть выше или ниже цены альтернативы, а разницу в стоимости можно оценить путем построения *модели обменного курса* [Smith, 2012]. Модель обменного курса оценивает потенциальную разницу в цене между двумя продуктами путем анализа и оценки отдельных характеристик продукта, которые могут быть выгодными (способствовать повышению цены) или невыгодными (способствовать снижению цены), как показано на рис. 6.2. Окончательную цену, также известную как курс обмена, можно получить путем сложения этой дифференциальной стоимости и базовой цены альтернативы.

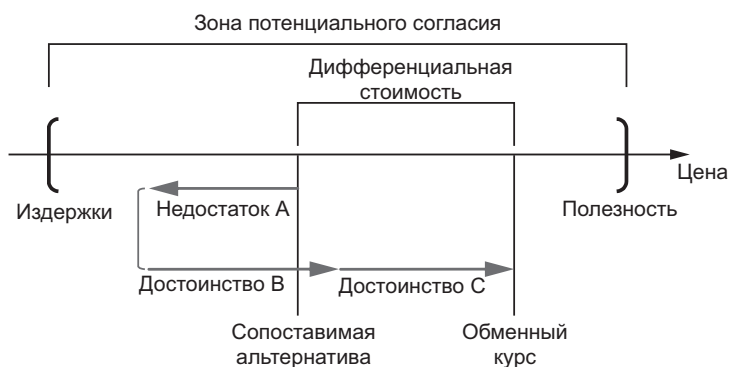


Рис. 6.2. Границы цены и курс обмена

Точный дизайн модели обменного курса зависит от характера продуктов и их различий. Во многих случаях характеристики продуктов можно оценить с помощью таких методологий, как *совместный анализ* (conjoint analysis), опирающийся на опрос потребителей [Green and Srinivasan, 1978]. В некоторых случаях модель можно создать путем анализа возможных результатов выбора того или иного продукта. Рассмотрим два примера:

- Дифференциальную стоимость более надежного продукта относительно менее надежного варианта можно оценить, приняв во внимание вероятность отказа (failure) и потенциальную стоимость замены. Если цена альтернативы p_A и стоимость замены p_R известны, тогда цену нового надежного продукта можно приблизительно оценить следующим образом:

$$p = p_R \cdot (1 - \Pr(\text{failure})) + (p_A + p_R) \cdot \Pr(\text{failure}). \quad (6.5)$$

- В модель обменного курса можно также включить цену вспомогательных и дополняющих продуктов. Например, производители бритв и лезвий обычно

разрабатывают свою продукцию так, чтобы лезвия для бритвы одного бренда не подходили к бритвам другого бренда, и, следовательно, обменный курс лезвия повышается из-за относительно высокой стоимости перехода на использование бритвы другого бренда.

Вопросы, касающиеся обменного курса, можно учесть непосредственно в структуре цен, о чем мы подробно поговорим в следующих разделах.

6.3.2. Субъективная ценность

Понятие полезности может предполагать, что покупатели принимают решения, сопоставляя свою готовность платить с предложенной ценой: продукт приобретается тогда и только тогда, когда цена ниже полезности. Однако такое «рациональное поведение» не является адекватной моделью реальных потребителей. Оценка ценности — это субъективный процесс, во многом зависящий от того, как именно ценность и цена сообщаются потенциальному покупателю и как тот воспринимает их. Неправильная подача ценности или цены может сформировать неправильные ожидания и сместить границы цены в нежелательном направлении. Эффективная подача, напротив, может повысить субъективную ценность продукта или уменьшить ценность сопоставимых альтернатив.

Подача ценности, на первый взгляд, может показаться не очень актуальной задачей в контексте алгоритмических методов, поскольку имеет дело с психологическими аспектами восприятия стоимости и цены, которые едва ли возможно закодировать в программном обеспечении. Однако оказывается, что анализ таких психологических моделей может дать практические правила для включения в ценовые структуры и, следовательно, учесть этот аспект в задачах оптимизации цен.

Одной из наиболее прочных основ, охватывающей многие важные аспекты восприятия стоимости и цены, является теория перспектив [Kahneman and Tversky, 1979]. Теория перспектив рассматривает процесс оценки с точки зрения оценки риска и может быть охарактеризована следующими положениями.

ТОЧКА ОТСЧЕТА. Потенциальные выгоды и потери, связанные со сделкой, оцениваются относительно некоторой точки отсчета. Точка отсчета основана на прошлом опыте (например, последняя наблюдаемая цена для данного или аналогичного продукта) или суждении и имеет тенденцию быть постоянной после ее выбора.

СНИЖЕНИЕ ЧУВСТВИТЕЛЬНОСТИ. Изменения выгод или потерь остро воспринимаются в зоне вокруг точки отсчета, но становятся менее заметными по мере их увеличения. Разница между скидками в 9 и 19 долларов представляется

существенной, но та же разница в десять долларов не воспринимается столь же значимой для скидок в 719 и 729 долларов.

НЕПРИЯТИЕ ПОТЕРЬ. Потери, как правило, воспринимаются острее, чем выгоды той же величины. Потеря 100 долларов, как правило, воспринимается как нечто более существенное, чем выгода в 100 долларов.

НЕПРИЯТИЕ РИСКА ДЛЯ ПОЛУЧЕНИЯ ВЫГОДЫ. Гарантированная выгода предпочтительнее конъюнктурной выгоды той же величины. Потенциальный клиент, стоящий перед выбором гарантированно получить 450 долларов или выиграть 1000 долларов с 50 % вероятностью (и, соответственно, с 50 % вероятностью не выиграть ничего), обычно отдает предпочтение первому варианту.

ПОЛОЖИТЕЛЬНОЕ ОТНОШЕНИЕ К РИСКУ ПОТЕРЬ. В отличие от выгод, потенциальные потери предпочтительнее гарантированных. Потенциальный клиент, стоящий перед выбором гарантированно потерять 450 долларов или потерять 1000 долларов с 50 % вероятностью (и, соответственно, с 50 % вероятностью не потерять ничего), обычно отдает предпочтение второй альтернативе.

Положения, изложенные выше, предполагают определенную форму зависимости между реальными и предполагаемыми выгодами и потерями, как показано на рис. 6.3. Наклон кривой в отрицательной зоне круче, чем в положительной, согласно гипотезе неприятия потерь, но на обоих концах крутизна уменьшается соответственно принципу снижения чувствительности.

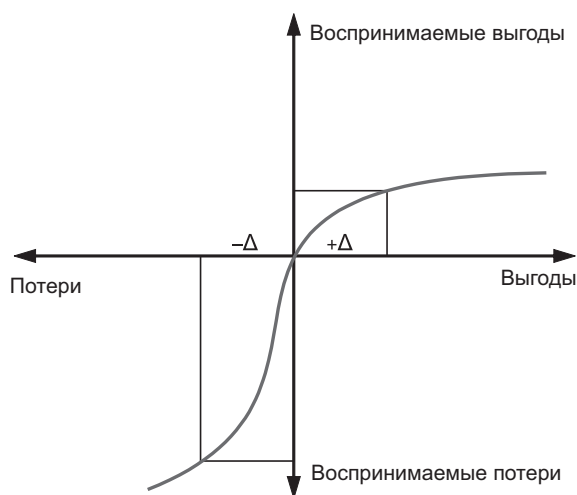


Рис. 6.3. Функция ценности в теории перспектив. Фактическое приращение выгоды $+\Delta$ может восприниматься как маленькая потеря, тогда как уменьшение на ту же величину $-\Delta$ может восприниматься как большая потеря

Теория перспектив предлагает несколько важных рекомендаций, которые можно использовать для оптимизации структуры цен:

- Повышение цен (надбавка) рассматривается гораздо более негативно, чем падение (скидки). По этой причине цены обычно сообщаются в виде объявленной цены и скидки. Это позволяет поддерживать постоянную цену по прейскуранту и манипулировать величиной скидки, в том числе для персонализированного ценообразования, без использования явных надбавок.
- Точка отсчета должна поддерживаться на высоком уровне. Это дополняет предыдущий пункт об объявленных ценах и скидках, поскольку снижение базовых цен может привести к нежелательному сдвигу точки отсчета.
- В общем случае лучше разделить выгоду на множество маленьких выгод. Одна большая выгода обесценивается из-за эффекта снижения чувствительности, поэтому структура цен с несколькими маленькими выгодами может иметь более высокую воспринимаемую ценность.
- Положительное отношение к риску потерь предполагает, что некоторые явные преимущества продукта можно заменить эквивалентными потенциальными выгодами без существенной потери воспринимаемой ценности. Например, мебель из плоской упаковки требует времени и усилий на ее сборку, но снижение цены и стоимости доставки все равно будет восприниматься как безусловные выгоды.

Мы используем некоторые элементы теории перспектив в оптимизации цен, чтобы учесть поведенческие факторы, несовместимые с базовым принципом рационального потребителя.

6.4. Цена и спрос

Полезность, модель курса обмена и другие методы оценки помогают предсказать ожидаемую готовность платить за данный продукт или услугу. Также можно вспомнить, что теория потребительского выбора (см. раздел 2.6.1) предлагает инструменты для прогнозирования потребительского выбора в случае нескольких доступных альтернатив. Конечно, можно попытаться продать продукт каждому клиенту по цене, напрямую вытекающей из индивидуальной готовности платить, но все-таки начнем с традиционной задачи определения общей структуры цен для всех клиентов. Это означает, что нам придется иметь дело с тысячами или даже миллионами индивидуальных решений о покупке, которые можно описать в вероятностных терминах.

Определим *готовность платить* за данный товар или услугу как максимальную цену, приемлемую для клиента. Клиент будет покупать продукт тогда и только

тогда, когда цена ниже той суммы, которую он готов заплатить. Популяцию клиентов можно описать, используя распределение готовности платить $w(p)$: для каждой пары цен, p_1 и p_2 , доля клиентов $f(p_1, p_2)$, готовых заплатить сумму между p_1 и p_2 , составляет

$$f(p_1, p_2) = \int_{p_1}^{p_2} w(p) dp. \quad (6.6)$$

Функцию спроса $q(p)$, также называемую *функцией цена/отклик*, можно выразить через $w(p)$ следующим образом:

$$q(p) = Q_{\max} \cdot \int_p^{\infty} w(x) dx, \quad (6.7)$$

где $Q_{\max} = q(0)$ — максимально достижимый спрос для данного продавца. Функцию спроса можно рассматривать не только как совокупную рыночную метрику, определяемую дисперсией готовности платить, но и как модель поведения одного клиента, в том смысле, что данный потребитель может быть готов купить разное количество продукта в зависимости от цены. В последнем случае готовность платить следует рассматривать как готовность платить *за единицу*.

Математический анализ функции спроса может дать нам дополнительную информацию и помочь определить полезные метрики и свойства. Во-первых, рассмотрим подробнее производную функции спроса:

$$\frac{\partial}{\partial p} q(p) = -Q_{\max} w(p). \quad (6.8)$$

Так как $w(p)$ неотрицательная, производная не является положительной для любого значения p , то есть функция спроса имеет наклон вниз. Наклон функции спроса, заданный ее производной, является мерой чувствительности к цене. Крутой уклон означает, что клиенты чувствительны к изменениям цены и спрос быстро падает с ее ростом, тогда как пологий наклон означает, что клиенты относительно нечувствительны к изменениям. Однако чувствительность к цене чаще измеряют не наклоном функции спроса, а *эластичностью спроса*, определяемой как отношение процентного изменения спроса к процентному изменению цены:

$$\varepsilon = -\frac{\Delta q / q}{\Delta p / p} = -\frac{p}{q(p)} \times \frac{\partial}{\partial p} q(p). \quad (6.9)$$

Согласно уравнению 6.9, эластичность является функцией цены и может быть разной в разных точках кривой спроса, однако этот термин часто используется более свободно, в предположении, что эластичность примерно постоянна в интересующем диапазоне, поэтому спрос можно охарактеризовать одним значением ε . Эластичность спроса не зависит от величины цены или ценности, поэтому

обеспечивает удобный способ измерения и сравнения чувствительности к ценам. Эластичные рынки, с $\varepsilon > 1$, реагируют на небольшое изменение цены значительным изменением спроса. Например, согласно исследованиям, эластичность ресторанных блюд составляет около 2,3, то есть увеличение цены на 10 % может вызвать снижение спроса на 23 %. Неэластичные рынки, с $\varepsilon < 1$, реагируют на изменения цен небольшим изменением спроса. Например, эластичность цен на автомобильный бензин в США оценивалась примерно в 0,04, то есть чтобы уменьшить количество автомобильных поездок на 1 %, требуется увеличить цену на бензин на 50 %. Однако следует различать эластичность для категорий товаров и для отдельных брендов в категории. Как правило, одну категорию трудно заменить другой категорией, поэтому во многих сферах спрос на уровне категории является относительно неэластичным. Смена бренда происходит проще, что делает кривые спроса более эластичными с точки зрения одного продавца.

Теперь рассмотрим несколько часто используемых моделей спроса, выразив их в терминах $w(p)$ и $q(p)$.

6.4.1. Линейная кривая спроса

Простую модель спроса можно получить, предположив, что готовность платить равномерно распределена в диапазоне от 0 до максимально приемлемой цены P :

$$w(p) = \text{unif}(0, P) = \begin{cases} 1/P, & 0 \leq p \leq P \\ 0 & \text{в противном случае.} \end{cases} \quad (6.10)$$

Функцию спроса можно получить интегрированием $w(p)$, согласно уравнению 6.7:

$$\begin{aligned} q(p) &= Q_{\max} \int_p^P w(x) dx = \\ &= Q_{\max} \left(1 - \frac{p}{P} \right) = \\ &= -\frac{Q_{\max}}{P} \cdot p + Q_{\max}. \end{aligned} \quad (6.11)$$

Как видите, равномерно распределенная готовность платить приводит к линейной функции спроса, как показано на рис. 6.4. Для оптимизации базовых ценовых структур будем исходить из предположения о линейности кривых спроса из-за их аналитического удобства, хотя это, как правило, очень грубая аппроксимация функций реального спроса. Одним из недостатков линейной модели спроса является предположение, что каждый доллар изменения цены дает одно и то же увеличение спроса. Это, как правило, неверно, потому что чувствительность к ценам

обычно высока вблизи точки отсчета (конкурентоспособной цены) и уменьшается при удалении от нее.

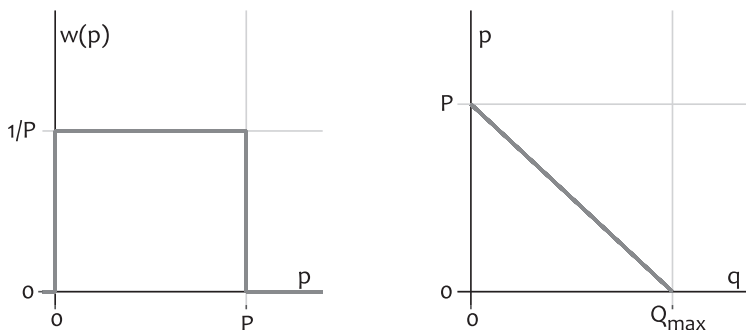


Рис. 6.4. Равномерное распределение готовности платить и соответствующая линейная кривая спроса. Обратите внимание: разместив цену на вертикальной оси на правом графике, мы следуем традиционной нотации, принятой в экономике, даже при том что спрос рассматривается как функция цены

6.4.2. Кривая спроса с постоянной эластичностью

Функцию спроса с постоянной эластичностью можно получить из определения эластичности, предположив, что эластичность глобально постоянна. Это означает, что мы должны решить следующее уравнение для $q(p)$:

$$\frac{p}{q(p)} \cdot \frac{\partial}{\partial p} q(p) = -\varepsilon. \quad (6.12)$$

Это дифференциальное уравнение, и его решением является семейство функций, заданных уравнением

$$q(p) = C \cdot p^{-\varepsilon}, \quad (6.13)$$

где $C > 0$ — произвольный коэффициент. Этот коэффициент, по сути, является параметром функции спроса, который выбирается так, чтобы график функции соответствовал известным точкам данных, а именно наблюдаемым парам цена/спрос. Рассчитать готовность платить, соответствующую спросу с постоянной эластичностью, можно, подставив 6.13 в 6.8:

$$w(p) = -\frac{\partial}{\partial p} q(p) \cdot \frac{1}{q(0)} = \varepsilon \cdot p^{-\varepsilon-1}. \quad (6.14)$$

Кривые спроса $q(p)$ и готовности платить $w(p)$ с постоянной эластичностью спроса изображены на рис. 6.5. Подобно линейной функции спроса, спрос с постоянной эластичностью может служить разумным приближением для относительно небольших изменений цен. Спрос с постоянной эластичностью правильно отражает плавное снижение готовности платить с ростом цены, но это также означает, что готовность платить — напомним, что это *максимально* приемлемая цена, — сосредоточена вблизи нуля, что не обязательно является реалистичным предположением.

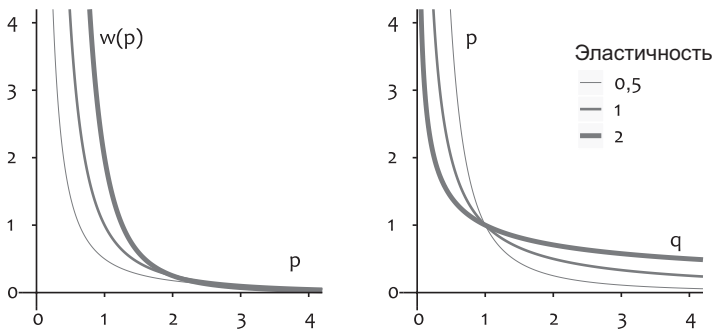


Рис. 6.5. Функция спроса с постоянной эластичностью $q(p) = p^{-\epsilon}$ и соответствующая готовность платить $w(p) = \epsilon \cdot p^{-\epsilon-1}$ для разных значений ϵ

6.4.3. Логит-кривая спроса

Логит-функция спроса пытается преодолеть ограничения линейной модели и модели с постоянной эластичностью, принимая во внимание тот факт, что эластичность цен достигает максимума вблизи точки отсчета. В частности, можно ожидать, что спрос на продукт будет оставаться устойчиво низким, если его цена значительно превышает цены конкурентов, и незначительные изменения в цене будут давать очень небольшой стимул, то есть локальная эластичность спроса будет низкой. Аналогично, если цены установлены намного ниже рыночных, это, скорее всего, стимулирует устойчиво высокий спрос, относительно нечувствительный к небольшим изменениям цен — клиенты все равно будут воспринимать покупку выгодной. Наиболее высокая чувствительность к ценам, вероятно, будет находиться в зоне вокруг цен конкурентов, где небольшие изменения цены могут существенно повлиять на спрос. Эти соображения, а также эмпирические наблюдения, предполагают сигмоидную форму кривой спроса, как показано на рис. 6.6. Сигмоидную кривую спроса, также называемую логит-функцией спроса (logit demand function), можно задать следующим образом:

$$q(p) = Q_{\max} \cdot \frac{1}{1 + e^{a+bp}}, \quad (6.15)$$

где Q_{\max} — максимально достижимый спрос, а b — параметр, управляющий крутизной кривой спроса. Для любых a и b максимальная чувствительность к ценам будет достигнута при цене, численно равной $-(a/b)$, поэтому параметр a можно использовать для сдвига точки отсчета, если задан b . Параметры Q_{\max} , a и b можно оценить, аппроксимировав логистическую кривую по наблюдаемым точкам данных.

Функцию готовности платить для логит-функции спроса легко получить путем дифференциации спроса:

$$w(p) = -\frac{\partial}{\partial p} q(p) \cdot \frac{1}{q(0)} = b(1 + e^a) \frac{e^{a+bp}}{(1 + e^{a+bp})^2}. \quad (6.16)$$

Как показано на рис. 6.6, логистическая готовность платить — это колоколообразная кривая, похожая на кривую нормального распределения.

Логит-функция спроса тесно связана с полиномиальной логит-моделью (Multinomial Logit, MNL), которую мы рассмотрели в разделе 2.6.1.1. Напомню, что если данный клиент n выбирает продукт или услугу из нескольких альтернатив $(1, \dots, J)$ и полезность от выбора варианта i измеряется как V_{ni} , тогда модель MNL утверждает, что выбор варианта i имеет следующую вероятность:

$$P_{ni} = \frac{e^{V_{ni}}}{\sum_{j=1}^J e^{V_{nj}}}. \quad (6.17)$$

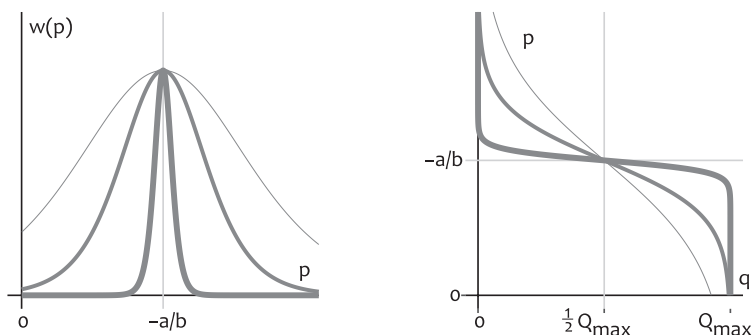


Рис. 6.6. Логит-функция спроса для нескольких значений параметра b и соответствующая кривая готовности платить. Толщина линий пропорциональна величине b

Полезность V_{ni} обычно измеряется с помощью регрессионной модели, учитывающей различные свойства клиентов и продуктов. Однако если моделировать ее

как линейную функцию цены $V_{mi} = -b_j p_j$, где b_j — коэффициенты регрессии, тогда вероятность P_{mi} не будет зависеть от индивидуальных потребителей и станет равной средней вероятности выбора продукта — рыночной доли продукта μ_i :

$$\mu_i(p) = \frac{e^{-b_i p_i}}{\sum_j e^{-b_j p_j}}. \quad (6.18)$$

Уравнение 6.18 — это, по сути, новая модель спроса, еще более гибкая, чем наша логит-модель, потому что позволяет явно учитывать индивидуальные конкурентные цены. Однако если конкурентные цены рассматривать как единый параметр, который может быть определен как

$$c = \sum_{j \neq i} e^{-b_j p_j}, \quad (6.19)$$

то долю рынка данного продукта i можно выразить следующим образом (обратите внимание, что здесь используется тождество $c^{-1} = e^{-\ln c}$):

$$\begin{aligned} \mu_i(p_i) &= \frac{e^{-b_i p_i}}{\sum_{j \neq i} e^{-b_j p_j} + e^{-b_i p_i}} = \frac{e^{-b_i p_i}}{c + e^{-b_i p_i}} \cdot \frac{c^{-1}}{c^{-1}} = \\ &= \frac{e^{-\ln c - b_i p_i}}{1 + e^{-\ln c - b_i p_i}}, \end{aligned} \quad (6.20)$$

что то же самое, что и логит-модель спроса, определяемая уравнением 6.15.

6.5. Базовые структуры цен

Кривая спроса описывает взаимосвязь между ценой и объемом спроса. Это позволяет выразить зависимость доходов и прибыли фирмы от цены и решить задачу оптимизации для определения оптимального уровня цены. Такой подход может выглядеть как точный способ расчета оптимальных цен, но он редко дает приемлемые результаты, потому что трудно, а порой невозможно точно оценить кривую спроса, которая учитывала бы все последствия изменения цен, включая реакцию конкурентов и другие стратегические решения. Однако формализация задачи оптимизации может подсказать полезные идеи и обеспечить поддержку в принятии решений, что является важным шагом на пути к программному решению. Анализ кривых спроса также помогает обосновать различные ценовые структуры и их ключевые свойства, что необходимо для реализации более продвинутых и автоматизированных оптимизаций, которые мы рассмотрим в последующих разделах.

6.5.1. Цена за единицу

Первая структура цен, которую мы рассмотрим, — это цены на отдельные элементы, или единицы, например на одну книгу, одну рубашку или один килограмм апельсинов. Сначала запишем стандартное уравнение прибыли фирмы на основе функции спроса $q(p)$:

$$G = q(p) \cdot (p - V), \quad (6.21)$$

где G — прибыль, p — цена, а V — переменные затраты. В целях упрощения я опустил фиксированные затраты. Напомню также, что линейная кривая спроса определяется уравнением

$$q(p) = Q_{\max} \cdot \left(1 - \frac{P}{p}\right), \quad (6.22)$$

где P — максимальная сумма, которую потребитель готов заплатить, и, следовательно, максимально приемлемая цена. Цену можно оптимизировать, взяв производную от прибыли относительно цены и приравняв ее к нулю:

$$\frac{\partial G}{\partial p} = \frac{\partial q}{\partial p}(p) \cdot (p - V) + q(p) = 0. \quad (6.23)$$

Решая это уравнение для p , получим оптимальную цену, равную среднему для P и V :

$$P_{opt} = \frac{P + V}{2}. \quad (6.24)$$

Оптимальную цену можно подставить в уравнение 6.22 и определить количество единиц, которые фирма, как ожидается, сможет продать по этой цене:

$$q_{opt} = \frac{Q_{\max}}{2P} (P - V). \quad (6.25)$$

Наконец, прибыль от выбора этой цены составит

$$G_{opt} = \frac{Q_{\max}}{4P} (P - V)^2. \quad (6.26)$$

Геометрическая интерпретация уравнений, приведенных выше, показана на рис. 6.7. Обратите внимание, что прибыль численно равна площади прямоугольника, ограниченного значениями p_{opt} и V .

Аналогичный результат можно получить для спроса с постоянной эластичностью. Согласно определению эластичности, получаем

$$\frac{\partial q}{\partial p}(p) = -\varepsilon \cdot \frac{q(p)}{p}. \quad (6.27)$$

Подставив это выражение в уравнение 6.23, найдем оптимальную цену:

$$p_{opt} = V \cdot \frac{\varepsilon}{\varepsilon - 1}. \quad (6.28)$$

Уравнение 6.28 очень удобно для демонстрации некоторых слабых сторон стратегической оптимизации цен с базовыми кривыми спроса. Рассмотрим пример фирмы, производящей товар стоимостью 10 долларов и желающей определить оптимальную цену продажи с использованием расчетной эластичности. Если предположить, что по данным продаж эластичность оценивается как 1,5, тогда оптимальная цена будет равна 30 долларам. Однако эластичность трудно оценить с высокой точностью, и вполне вероятно, что оценка на самом деле находится в диапазоне $1,5 \pm 0,4$. Это дает диапазону «оптимальных» цен от 21 до 110 долларов, как показано в табл. 6.2. Понятно, что этот результат имеет ограниченную практическую ценность.

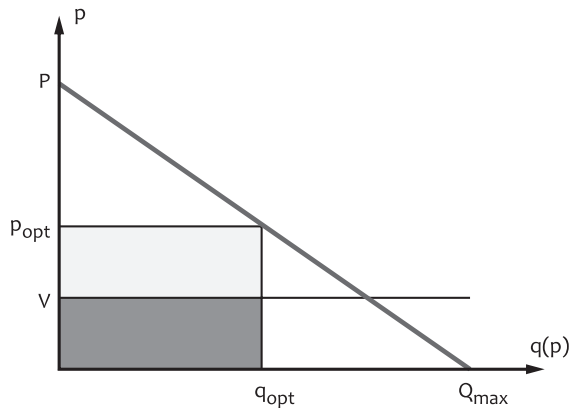


Рис. 6.7. Оптимизация цены за единицу для линейной кривой спроса

Таблица 6.2. Оптимальные цены для разных значений эластичности и переменных затрат в 10 долларов, рассчитанные с помощью модели эластичности затрат/спроса

Эластичность	1,10	1,20	1,30	1,40	1,50	1,60	1,70	1,80	1,90
Оптимальная цена	110	60	43	35	30	27	24	23	21

6.5.2. Сегментация рынка

Практически все рынки демонстрируют неоднородность спроса, обусловленную тем, что разные клиенты, и даже одни и те же клиенты в разные моменты времени, оценивают продукты по-разному. Такое разнообразие оценок объясняется многими причинами. Потребительские рынки, например, в принципе неоднородны, потому количество человеческих потребностей, таких как питание или одежда, относительно ограничено, но доходы сильно различаются, что приводит к очень разным суммам, потраченных разными людьми на одни и те же потребности. Клиенты могут использовать одни и те же или похожие продукты по-разному и извлекать разные ценности из свойств продукта, иметь более или менее подробную информацию о конкурентных предложениях и т. д. Мы уже видели, как эта неоднородность создает благодатное поле для целевого продвижения и рекламы, и теперь можем исследовать ее влияние на ценовые решения. К счастью, анализ оптимизации цены за единицу дает для этого очень удобную основу.

Как показано на рис. 6.7, максимально достижимый доход численно равен общей площади под кривой спроса, поэтому максимально достижимую прибыль можно оценить как

$$G_{\max} = \frac{1}{2} P \cdot Q_{\max}. \quad (6.29)$$

В то же время любая отдельно взятая цена p_{opt} , оптимальная или неоптимальная, является компромиссом, потому что некоторые клиенты не будут покупать продукт, если посчитают его слишком дорогим, но готовы будут купить его по более низкой цене, между p_{opt} и V , и тем самым положительно повлиять на прибыль. С другой стороны, некоторые клиенты готовы платить цену выше p_{opt} , однако генерируемый ими объем продаж будет относительно невелик. В обоих случаях фирма не получит дополнительных прибылей, которые лежат в треугольнике между кривой спроса и линией переменных издержек. Ценовая сегментация — это естественный способ преодоления ограничений единой цены путем сегментирования клиентов по суммам, которые они готовы платить, и предложения разных цен разным сегментам. На рис. 6.8 показан частный случай этой стратегии, когда обычная цена была дополнена более высокой премиальной ценой и более низкой дисконтной ценой. Обратите внимание, как увеличивается область прибыли в сравнении со стратегией одной цены.

Это соображение ведет нас к сложному вопросу: как продавать одни и те же или аналогичные продукты разным клиентам по разным ценам? В широком смысле это требует установки барьеров между клиентами с разной готовностью платить, чтобы клиенты с большей готовностью не смогли платить более низкую цену, пред-

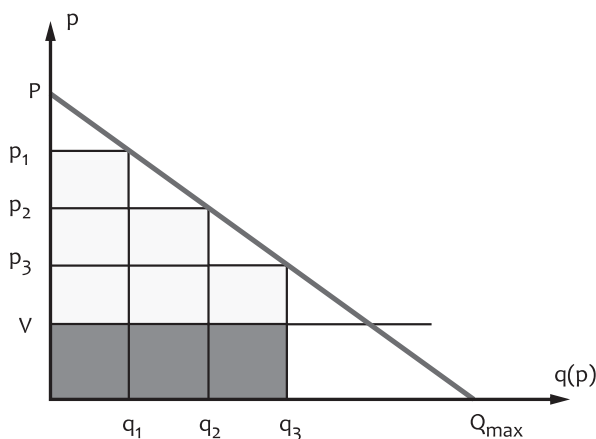


Рис. 6.8. Оптимизация прибыли путем сегментирования цен

назначенную для сегментов с меньшей готовностью. Существует много механизмов ограждения, и в разных сферах они разные, хотя большинство вариантов можно свести к нескольким основным принципам. Рассмотрим несколько примеров ценового ограждения из розничной торговли, демонстрирующих замечательную изобретательность в этом отношении.

ТОРГОВЫЕ ЗОНЫ. Магазины розничной торговли, как правило, расположены в разных районах с различными демографическими и конкурентными факторами, такими как средний доход домохозяйства, средняя численность семьи, расстояние до ближайшего магазина-конкурента и т. д. Это, естественно, разделяет клиентов по уровню ценовой чувствительности и способности или готовности искать альтернативного продавца, что позволяет продавцу устанавливать цены на уровне магазина, отличающиеся в разных зонах.

РАЗМЕР УПАКОВКИ. Ходовые потребительские товары (Fast-Moving Consumer Goods, FMCG), такие как безалкогольные напитки или туалетные принадлежности, имеют высокие показатели оборота, и потребители, естественно, могут выбирать между частой покупкой небольших количеств продукта и редким приобретением больших объемов. На этот компромисс также влияют демографические факторы, такие как размер домохозяйства. Это создает барьеры по признаку готовности покупать большие или малые упаковки и устанавливает различные допустимые пределы для упаковок разного размера. Предложения и скидки, основанные на количестве и частоте покупок, также относятся к этой категории.

РАСПРОДАЖИ. Клиентов можно разделить по их готовности ждать более низкой цены и готовности купить продукт сразу, но по обычной цене. Этот тип сегмента-

ции широко используется в торговле одеждой, где сезонные распродажи являются одним из основных механизмов маркетинга.

КУПОНЫ. Многие клиенты порой не готовы купить данный продукт по обычной цене, но могут задуматься о его приобретении по сниженной цене. В результате продавец может извлечь дополнительную выгоду от скидок, поскольку они привлекают дополнительных клиентов, даже притом что маржа в этом случае ниже, по сравнению с покупательскими привычками постоянных клиентов. С другой стороны, предлагать скидку чрезмерно широкой аудитории может быть вредно, так как ею будут пользоваться даже те клиенты, которые готовы платить обычную цену (без скидки). Методы моделирования отклика, описанные в главе 3, помогают решить эту проблему. Однако существует традиционное решение, которое используется с XIX века, — купоны. Купон — это скидка, для получения которой требуется приложить определенные усилия (например, клиент должен найти купон в газете, вырезать и представить в магазине), что отделяет клиентов с готовностью тратить время и усилия на получение скидки.

КАНАЛЫ ПРОДАЖ. Каналы продаж представляют естественные ограждения, потому что клиенты выбирают каналы по критериям, сильно коррелирующим с их готовностью платить. Например, ценовая чувствительность покупателей винных магазинов неизменно ниже, чем у покупателей, приобретающих то же вино в продуктовых магазинах [Cuellar and Brunamonti, 2013].

ОТДЕЛЫ. Продавцы и производители часто дифференцируют наценки в соответствии с различиями в ценовой чувствительности между полом и возрастом. Например, женская одежда, как правило, дороже мужской.

КЛУБНЫЕ КАРТЫ. Членство помогает отличить случайных покупателей от высокодоходных постоянных клиентов, для которых непубличность покупок ценнее членских взносов.

БРЕНДИНГ. Розничные торговцы и производители создают отдельные бренды, ориентированные на потребительские сегменты с более высокой или низкой готовностью платить по отношению к основному бренду. Бренды могут позиционироваться как менее престижные, чтобы продавать продукцию дешевле и не создавать конкуренцию основному бренду, или как более престижные, чтобы получить дополнительный доход от потребителей с более низкой ценовой чувствительностью.

Этот список ценовых барьеров можно дополнить другими методами из других отраслей, такими как тарифные классы в авиакомпаниях или сделки с кредитными картами. Стратегию оптимизации ценового сегментирования можно продемонстрировать примером разбивки цен на n сегментов с целью максими-

зирать прибыль. Начнем со следующего уравнения, которое прямо следует из рис. 6.8:

$$G = \sum_{i=1}^n (q_i - q_{i-1})(p_i - V), \quad (6.30)$$

где p_i и q_i — цена и количество проданного товара для сегмента i соответственно, и $q_0 = 0$. Количество проданного товара по цене p_i составит:

$$q_i = Q_{\max} \left(1 - \frac{p_i}{P} \right). \quad (6.31)$$

Мы можем найти цены, максимизирующие прибыль, взяв частные производные от G и приравняв их к нулю. Подставив уравнение 6.31 в уравнение 6.30, установив $p_0 = S$ и $p_{n+1} = V$ и выполнив алгебраические упрощения, найдем

$$\frac{\partial G}{\partial p_i} = \frac{Q_{\max}}{P} (p_{i-1} - 2p_i + p_{i+1}), \quad 1 \leq i \leq n. \quad (6.32)$$

Приравнявая частные производные к нулю, находим рекуррентную зависимость для цен сегмента:

$$p_i = \frac{p_{i-1} + p_{i+1}}{2}. \quad (6.33)$$

Легко проверить, что это соотношение, а также начальные условия $p_0 = S$ и $p_{n+1} = V$ удовлетворяются следующими ценами сегмента:

$$p_i^{opt} = \frac{1}{n+1} [(n+1-i) \cdot P + i \cdot V]. \quad (6.34)$$

Следовательно, оптимальные цены должны равномерно распределяться между переменными затратами V и максимально приемлемой ценой P . Это иллюстрирует пример на рис. 6.9.

Определить полученную прибыль с этими ценами можно, подставив уравнение 6.34 в уравнение 6.30:

$$G_{opt} = Q_{\max} \cdot \frac{n(P-V)^2}{2(n+1)P}. \quad (6.35)$$

Это уравнение является обобщением уравнения прибыли для цены за единицу продукции 6.26. Как видите, прибыль растет пропорционально $n/(n+1)$ по мере увеличения количества сегментов и приближается к максимально достижимой.

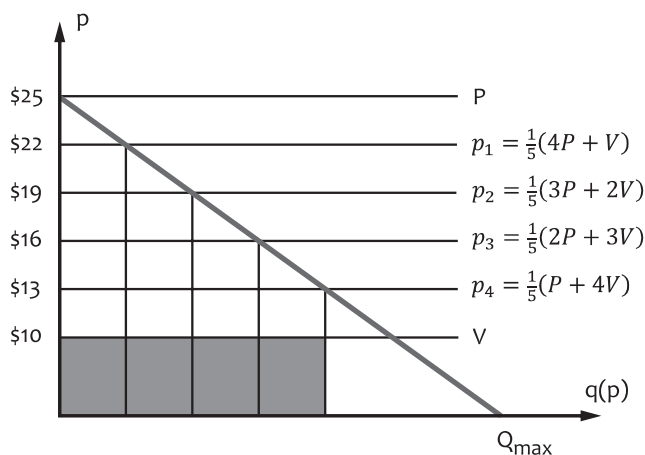


Рис. 6.9. Пример оптимальных цен для четырех сегментов с переменными затратами 10 долларов и максимально допустимой ценой 25 долларов

Ценовая сегментация является, пожалуй, наиболее мощным и широко используемым методом ценообразования. Однако его эффективность всецело зависит от способности разграничивать сегменты. Строгая сегментация достижима в относительно редких случаях. Например, парк аттракционов может устанавливать разные цены на билеты для разных возрастных групп, проверяя возраст посетителей по документам.

Однако большинство решений сегментирования далеки от совершенства и дают возможность клиентам из более высокого сегмента купить продукт, продаваемый клиентам из более низкого сегмента по более низкой цене. Например, клиенты онлайн-магазинов, как известно, имеют более высокую ценовую чувствительность, чем клиенты обычных магазинов, потому что имеют более простую возможность сравнивать цены и другие факторы, поэтому розничные торговцы довольно часто снижают цены в онлайн-магазинах. Это вызывает так называемое *поведение демонстрационного зала*, когда клиенты рассматривают объекты в магазине, а затем покупают их онлайн. Такая *каннибализация рынка*¹ может иметь весьма пагубные последствия, как показано в примере, приведенном в табл. 6.3. Добавление низкоценового сегмента увеличивает общую прибыль, когда сегментация проводится безупречно, но утечка 500 клиентов в этот сегмент из более высокого полностью съедает прибыль.

¹ Ситуация на рынке, при которой новый товар предприятия забирает покупателей старого товара этого же предприятия, то есть «съедает» собственную долю рынка. — *Примеч. ред.*

Таблица 6.3. Пример сегментации с каннибализацией спроса и без нее. Переменные затраты V составляют 10 долларов за штуку, и $\text{Выгода} = \text{Доход} - V \times \text{Спрос}$.

		Цена	Спрос	Доход	Выгода	Общая выгода
С одним сегментом		19	5000	95 000	45 000	45 000
Безупречная сегментация	A	19	5000	95 000	45 000	48 000
	B	13	1000	13 000	3000	
Ошибочная сегментация	A	19	4500	85 500	40 500	45 000
	B	13	1500	19 500	4500	

6.5.3. Комплексное ценообразование

Как мы убедились, сегментация является мощным методом извлечения дополнительных доходов из неоднородного спроса. Одним из популярных методов создания барьеров между клиентскими сегментами является использование различий в интенсивности и закономерностях использования продуктов. Например, производители фотокамер, как правило, предлагают широкий спектр продуктов, включая модели начального уровня, продвинутое камеры для энтузиастов и профессиональные камеры. Эти категории фотокамер существенно различаются по качеству получаемых изображений и другим функциональным свойствам, однако производители также пытаются получить прибыль от более интенсивного использования профессионального оборудования, предлагая повышенную долговечность и более широкие возможности. В связи с этим наиболее точной сегментации можно достичь, привлекая владельцев камер количеством снимков, которые они смогут получить, чтобы добиться прямой связи доходности с использованием. В отношении фотокамер этот подход может оказаться неосуществимым по причинам реализации и конкуренции, но его разновидности успешно применяются в других отраслях. Две структурно сходные ценовые стратегии, которые используют эту идею, — это *двухкомпонентные тарифы* и *связывающие соглашения*.

ДВУХКОМПОНЕНТНЫЕ ТАРИФЫ. Двухкомпонентный тариф — это ценовая структура с двумя компонентами — *входной платой* и *мерной ценой*. Входная плата взимается за доступ к продукту или услуге. Мерная цена — это плата за единицу, зависящая от потребленного количества. Классическим примером двухкомпонентных тарифов являются телекоммуникационные услуги, когда помимо платы за каждую использованную минуту или гигабайт взимается плата за подключение. Другими примечательными примерами могут служить коммунальные услуги, такие как электроснабжение, газоснабжение или водоснабжение, корпоративное

ПО, к базовой цене которого часто добавляется сумма, пропорциональная числу пользователей, заказов или одновременных соединений, и парки развлечений, которые могут взимать плату за вход и цену за пользование аттракционами.

СВЯЗЫВАЮЩИЕ СОГЛАШЕНИЯ. Некоторые продукты очень тесно связаны друг с другом, в том смысле что клиент не может извлечь большой ценности из одного продукта без другого. Это позволяет производителю создавать механизмы связывания, не позволяющие клиенту переключаться между брендами и балансировать цены на соответствующие продукты. Примеры связывающих соглашений часто можно найти, когда долговечное изделие (связывающий продукт) дополняется расходной частью (связанный продукт), как бритвенный станок и бритвенные лезвия или картриджи для принтера и чернила. Доходы, получаемые от расходных частей, могут занимать доминирующее положение в пожизненной ценности потребителя, поэтому долговечный продукт может стоить меньше или даже меньше его себестоимости.

Теперь рассмотрим количественные модели, иллюстрирующие оптимизацию входной платы p_e и мерной цены p_m в двухкомпонентном тарифе. Выше упоминалось, что кривую спроса можно интерпретировать как совокупный спрос, определяемый максимальной готовностью платить и уровнем потребления одного клиента при заданной цене. Двухкомпонентные тарифы прямо зависят от уровня потребления и, следовательно, требуют учета не только уровня потребления, но и неоднородности спроса, поэтому мы должны использовать более сложную модель, чем та, что использовалась для удельного ценообразования и ценовой сегментации [Smith, 2012; Oi, 1971].

Для начала рассмотрим случай одного потребителя, кривая спроса для которого изображена на рис. 6.10.

Напомним, что уравнение линейной кривой спроса имеет вид

$$q = Q_{\max} \left(1 - \frac{P_m}{P} \right). \quad (6.36)$$

Потребитель оценивает продукт или услугу по цене P , поэтому профицит от покупки одной единицы по мерной цене p_m численно равен площади треугольника под кривой спроса, ограниченного линией p_m . Следовательно, можно предположить, что оптимальная входная плата равна этому профициту, потому что более низкая плата приведет к недополучению доступной прибыли, а более высокая — вытеснит потребителя с рынка. Это означает, что плату можно выразить как площадь под кривой спроса:

$$p_e^{opt} = \frac{q}{2} (P - p_m) = \frac{Q_{\max}}{2P} (P - p_m)^2. \quad (6.37)$$

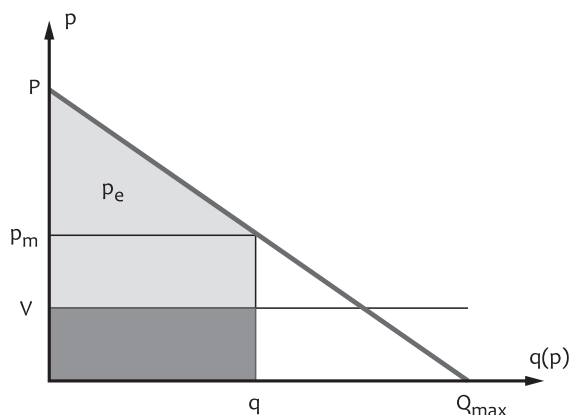


Рис. 6.10. Оптимизация цены для двухкомпонентного тарифа

Общая прибыль, полученная от клиента, определяется как сумма вступительной и мерной платы:

$$G = p_e^{opt} + q(p_m - V). \quad (6.38)$$

Оптимальную мерную цену можно вычислить, приравняв производную прибыли к нулю:

$$\frac{\partial G}{\partial p_m} = \frac{\partial p_e^{opt}}{\partial p_m} + \frac{\partial q}{\partial p_m}(p_m - V) + q = 0. \quad (6.39)$$

Подставляя уравнения 6.37 и 6.38 в уравнение 6.39 и решая его для p_m , находим, что оптимальная цена должна быть равна предельным затратам:

$$p_m^{opt} = V. \quad (6.40)$$

Таким образом, двухкомпонентное ценообразование побуждает продавца устанавливать вступительную плату как можно выше и снижать мерную цену до минимума, то есть прибыль будет извлекаться исключительно из вступительного взноса. Эта стратегия, например, широко используется парками развлечений, которые склонны взимать высокую плату за вход, а не за аттракционы.

Однако подход с высокой входной платой сталкивается с сильными сдерживающими факторами при наличии конкуренции или неоднородного спроса, когда клиенты готовы приобрести различные количества продукта по заданной цене. Эта ситуация изображена на рис. 6.11, где несколько кривых спроса имеют один и тот же наклон, но различаются количеством приобретенных товаров.

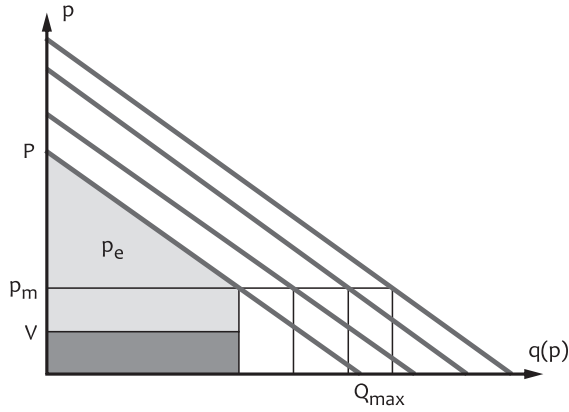


Рис. 6.11. Двухкомпонентный тариф в случае неоднородности потребления

Допустим, что каждая кривая спроса соответствует определенному потребительскому сегменту. Теперь мы не можем установить вступительный взнос выше профицита, соответствующего кривой самого низкого спроса, потому что иначе потеряем клиентов. Обозначим отношение спроса i -го сегмента к спросу низшего сегмента как k_i . Теперь уравнения для кривых спроса можно записать так:

$$q_i = Q_{\max} k_i \left(1 - \frac{P_m}{k_i P} \right), \quad k_i \geq 1. \quad (6.41)$$

Прибыль можно выразить как сумму входной платы с мерной ценой каждого из сегментов:

$$G = p_e + \sum_i \mu_i q_i (p_m - V), \quad (6.42)$$

где μ_i — доля сегмента i , то есть, если общее число клиентов равно N , тогда сегмент i содержит $N \cdot \mu_i$ клиентов. Оптимальную мерную цену можно найти, взяв производную прибыли по мерной цене и приравняв ее к нулю:

$$\frac{\partial G}{\partial p_m} = \frac{\partial p_e}{\partial p_m} + (p_m - V) \sum_i \mu_i \frac{\partial q_i}{\partial p_m} + \sum_i \mu_i q_i = 0. \quad (6.43)$$

Уравнение 6.37 для p_e все еще выполняется в предположении, что P — максимально приемлемая цена для кривой наименьшего спроса, поэтому ее можно вставить в уравнение выше и, используя тот факт, что сумма всех i равна единице, получить простое выражение для оптимальной мерной цены:

$$p_m^{opt} = V + P \sum_i \mu_i (k_i - 1) = V + P(\mathbb{E}[k] - 1), \quad (6.44)$$

где $\mathbb{E}[k]$ — взвешенное среднее множителей спроса сегмента. Этот результат показывает, что неоднородность уровней спроса приводит к росту мерной цены и, следовательно, снижает вступительную плату до уровня кривой самого низкого спроса, что может изменить баланс цен в исходном двухкомпонентном тарифе — мерный элемент может стать преобладающим над входным элементом в структуре общей прибыли. Этот сдвиг также может вызывать конкуренция, ограничивающая способность продавца извлекать потребительские излишки через вступительный взнос.

6.5.4. Пакетирование

В книгах по экономике под пакетированием часто понимается продажа двух или более отдельных продуктов в одной упаковке по общей цене. Это не очень точное определение, потому что практически любой продукт можно рассматривать как пакет из составляющих его частей, например, автомобиль — как пакет из двигателя, колес и других компонентов, которые можно продавать отдельно, по крайней мере, на промышленных рынках. С точки зрения оптимизации цен нас больше интересует *пакетное ценообразование*, когда два или более продуктов предлагаются по сниженной цене относительно суммы цен на отдельные продукты; данная скидка является единственным преимуществом перед альтернативой покупки разукрупненных товаров.

Пакетное ценообразование — популярная ценовая структура, которую можно найти во многих отраслях. В качестве примеров ценовых пакетов можно привести билеты на спортивные и театральные сезоны, ресторанные блюда, включающие закуску, основное блюдо и десерт, наборы с несколькими сумками разного размера, кухонные принадлежности и программные комплекты. Продавец, как правило, имеет выбор из трех вариантов: продавать продукты по отдельности, продавать группу продуктов только в виде пакета (этот вариант известен как *чистое пакетирование*) или предлагать продукты как по отдельности, так и в виде пакета (этот вариант называют *смешанным пакетированием*). Интуитивно понятно, что скидка за пакет должна уравниваться получением большей прибыли от продажи продуктов по отдельности. Выше мы видели, что сложные ценовые структуры часто достигают этого, используя неоднородность в готовности платить. Следовательно, можно сделать предположение, что, применяя подход пакетирования, можно использовать различия в готовности платить за различные товары.

ПРИМЕР 6.2

В качестве примера рассмотрим пакет офисного ПО, включающий приложение для работы с электронными таблицами и создания презентаций. Самый простой сценарий, который мы можем проанализировать, — это рынок с одним сегментом, в котором стоимость каждого приложения одинакова для всех клиентов. Предположим, что одна пользовательская лицензия на приложение для работы с электронными таблицами оценивается в 100 долларов, а на приложение для создания презентаций — в 150 долларов. В этом сценарии мы можем назначить цены для отдельных продуктов, а объединение двух продуктов в пакет не будет выгодным, потому что любая цена пакета, отличающаяся от 250 долларов, приведет к потерям. Сценарий с двумя сегментами, представленный в табл. 6.4, более интересен, потому что в разных сегментах готовность платить разная — в отделах продаж высоко ценят приложение для презентаций, а в бухгалтериях отдают предпочтение электронной таблице. Максимальные цены для отдельных продуктов, которые удержат оба сегмента на рынке, составляют 100 долларов за приложение электронной таблицы и 100 долларов за приложения подготовки презентаций. Если размеры двух сегментов равны, эти цены являются оптимальными и приносят прибыль в размере 400 долларов. Однако цену комплекта, включающего два продукта, можно установить на уровне 250 долларов, потому что оба сегмента, отделов продаж и бухгалтерий, оценивают пару продуктов суммой 250 долларов. Общая прибыль от двух сегментов в этом случае составит 500 долларов, что лучше 400, заработанных без пакетирования цен.

Таблица 6.4. Пример пакетного ценообразования с двумя продуктами и двумя сегментами потребителей

Сегмент клиентов	Готовность платить	
	электронная таблица	презентации
Отделы продаж	100	150
Бухгалтерии	150	100

В примере выше используется асимметрия готовности платить между двумя сегментами. Если все клиенты постоянно оценивают один товар выше другого, объединение продуктов в пакет не сможет дать более высокую прибыль, чем про-

даже этих продуктов по отдельности. Мы можем использовать это наблюдение для построения количественной модели оптимизации цены пакета. Преимущество этой модели заключается в том, что она делает относительно мало предположений о том, как распределяются клиентские сегменты, и может использовать преимущества численной оптимизации или моделирования для оптимизации ценообразования для произвольного числа сегментов или даже отдельных пользователей. Рассмотрим сценарий, где мы продаем два продукта, X и Y , и существует несколько клиентских сегментов, в которых максимальная готовность платить цену P_{ix} за продукт X равна или пропорциональна готовности платить цену P_{iy} за продукт Y , как показано на рис. 6.12. Для любых цен на продукты, p_x и p_y , потребительский сегмент попадает в одну из четырех областей (ничего не покупать, покупать только продукт X , покупать только продукт Y или покупать два продукта) в зависимости от соотношения между готовностью платить и соответствующей ценой. В нашем сценарии с положительно коррелированными оценками клиенты покупают либо все, либо ничего, поэтому пакетирование не добавляет ценности.

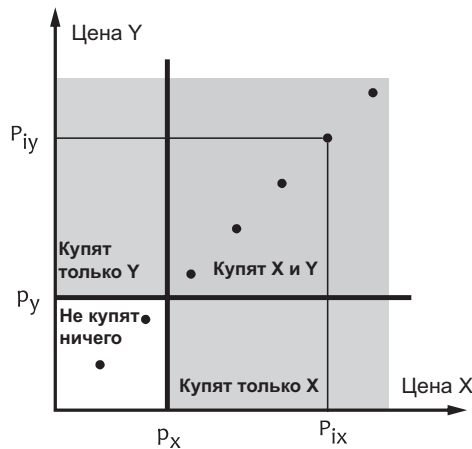


Рис. 6.12. Потребительские сегменты, в которых готовность платить за два продукта коррелирует положительно. Сегменты обозначены черными точками

Случай с асимметричной готовностью платить в корне отличается. В целях упрощения предположим, что все сегменты имеют одинаковую совокупную готовность платить, но сегменты могут распределять ее по продуктам по-разному, как показано на рис. 6.13. Именно так обстояло дело в примере с пакетом офисного программного обеспечения, который мы только что рассмотрели.

ЧИСТОЕ ПАКЕТИРОВАНИЕ. Цену пакета p_B можно установить равной совокупной готовности платить, так как она постоянна для всех сегментов:

$$p_{\beta} = p_{ix} + p_{iy}, \text{ постоянная для всех } i. \quad (6.45)$$

Прибыль от продажи пакета составит:

$$\begin{aligned} G_B &= \sum_i n_i (p_B - V_x - V_y) = \\ &= \sum_i n_i (p_{ix} + p_{iy} - V_x - V_y), \end{aligned} \quad (6.46)$$

где n_i — число клиентов в сегменте i , а V_x и V_y — переменные затраты для продуктов X и Y соответственно. С другой стороны, продажа продуктов по отдельности приносит следующую прибыль:

$$G_U = \sum_{i: p_{xi} \geq p_x} n_i (p_x - V_x) + \sum_{i: p_{yi} \geq p_y} n_i (p_y - V_y). \quad (6.47)$$

Первый и второй члены уравнения 6.47 представляют суммарную прибыль для продуктов X и Y соответственно. Сравнивая уравнения 6.46 и 6.47, нетрудно обнаружить, что G_B больше или равно G_U для любых цен на продукты p_x и p_y , поэтому в данном сценарии чистое пакетирование является более эффективной стратегией, чем продажа продуктов по отдельности.

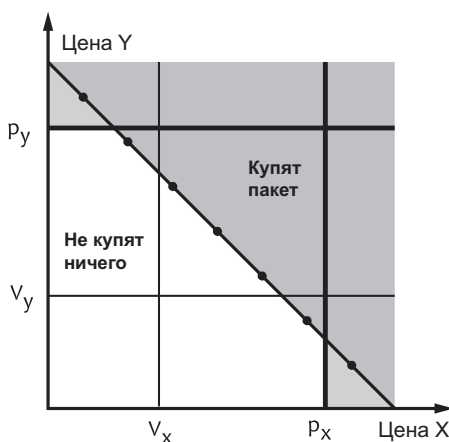


Рис. 6.13. Потребительские сегменты, в которых готовность платить за два продукта коррелирует отрицательно

СМЕШАННОЕ ПАКЕТИРОВАНИЕ. Можно попробовать расширить стратегию чистого пакетирования и организовать продажу продуктов как по отдельности, так и в пакете. Для этого нужно выбрать такие цены на отдельно продаваемые

продукты, чтобы не каннибализировать прибыль, получаемую от пакетирования. С этой целью можно сравнить прибыль от продажи пакета с прибылью от продажи отдельного продукта. Возьмем в качестве примера продукт X :

$$p_x - V_x > p_{\text{в}} - V_x - V_y. \quad (6.48)$$

Следовательно, цены на отдельные продукты должны быть

$$\begin{aligned} p_x &> p_{\text{в}} - V_y \\ p_y &> p_{\text{в}} - V_x. \end{aligned} \quad (6.49)$$

Это условие соответствует маленьким треугольникам на концах ценовой линии пакета на рисунке 6.13. Клиенты, попадающие в эти области, предпочтут вместо пакета купить либо продукт X (самый правый треугольник), либо продукт Y (самый левый треугольник), и будет генерировать дополнительный доход, более высокий, чем при чистом пакетировании.

Подход, описанный выше, обеспечивает относительно высокую гибкость для оптимизации пакетирования. Если предположить, что готовность платить оценивается для выборочных сегментов или отдельных потребителей, соответствующие точки можно закрепить на плоскости, а оптимальные цены найти с помощью численных методов оптимизации и разбиения плоскости на области с различной прибылью в зависимости от местоположений и размеров сегментов.

6.6. Прогнозирование спроса

Рассмотренные выше базовые модели спроса удобны для стратегического анализа, но часто слишком грубы для оптимизации фактических цен. Проблема в том, что спрос на тот или иной продукт зависит от многих факторов, в том числе собственных свойств продукта, таких как цена или бренд, цены конкурирующих продуктов в категории, распродажи и даже погода. Нам нужна более совершенная модель спроса, которая бы учитывала эти факторы и позволяла проводить анализ «что, если» для прогнозирования реакции на изменение цен, расширение и сокращение ассортимента, перераспределение площадей на полках. Эта модель является важным строительным блоком, который можно использовать во многих приложениях, которые зависят от количественной оценки спроса, в том числе следующие.

СТАТИЧЕСКАЯ ОПТИМИЗАЦИЯ ЦЕН. Для разных клиентских сегментов, определяемых каналами, местоположением и моделями предрасположенности, можно настроить свои базовые цены и уценки. Для этого в модель спроса нужно включить свойства сегмента.

ДИНАМИЧЕСКАЯ ОПТИМИЗАЦИЯ ЦЕН. Уценка, распродажа и ценообразование продуктов с ограниченными объемами поставок требуют оптимизации цен в зависимости от времени. То есть модель спроса должна также учитывать течение времени.

УПРАВЛЕНИЕ КАТЕГОРИЯМИ. Для оптимизации пространства на полках и ассортимента важно понимать зависимости спроса на разные продукты.

ОПТИМИЗАЦИЯ ЗАПАСОВ. Управление цепочками поставок и пополнением запасов выигрывает от моделирования спроса. Прогнозирование спроса особенно важно при планировании крупных и скоротечных распродаж.

Прогнозирование спроса можно рассматривать как относительно простую задачу анализа данных, которая сводится к построению регрессионной модели и ее применению к историческим данным. Однако разработать регрессионную модель не так просто, как может показаться, потому что на спрос влияют многие факторы со сложными зависимостями. Может потребоваться объединить несколько базовых функций спроса, таких как линейный спрос или логит-спрос, чтобы собрать достаточно гибкую модель, способную правильно фиксировать сезонные изменения, потребительский выбор, эластичность цен и другие факторы.

Проектирование моделей спроса — целое искусство, потому что разные задачи оптимизации требуют разных моделей прогнозирования спроса и вряд ли можно построить универсальную модель спроса, которая включает широкий спектр факторов, влияющих на спрос, таких как:

- Использование модели. Модель может включать или не включать элементы управления потребительским выбором, изменением спроса с течением времени, конкурентными ценами и т. д. Выбор элементов управления зависит от приложения, для которого создается модель.
- Доступные данные. Доступность, достоверность и полнота данных влияют на дизайн модели и ее возможности.
- Бизнес-область, модель и процесс. Модель спроса отражает терминологию, ограничения и структуру конкретного бизнеса. Например, модель спроса может предсказать норму спроса для отдельных продуктов или групп, включающих несколько вариантов продукта разных размеров, цветов или вкусов.
- Экспериментирование. Модели спроса, как и большинство реальных предиктивных моделей, явно и неявно включают много предметных знаний и требуют обширной настройки и проведения множества экспериментов.

Однако большое количество идей и приемов можно почерпнуть, изучая промышленные модели прогнозирования спроса. Структура моделей и выбор при-

знаков — это строительные блоки и подсказки, которые можно использовать для решения будущих задач прогнозирования спроса. Далее мы рассмотрим две реальные модели спроса из области розничной торговли, а также исследуем разницу между этими двумя примерами и несколькими моделями, применяемыми другими компаниями. Все эти модели были созданы в контексте оптимизации цен и ассортимента, поэтому они хорошо согласуются с методами оптимизации, которые мы рассмотрим далее в этой главе.

6.6.1. Модель спроса для оптимизации ассортимента

Первой мы рассмотрим модель спроса, разработанную для оптимизации ассортимента в сети супермаркетов Albert Heijn в Нидерландах [Kök and Fisher, 2007]. В ней особое внимание уделяется потребительским решениям, позволяющим провести детальный анализ факторов, влияющих на потребительский выбор.

Сеть супермаркетов продает большое количество продуктов, которые делятся на товарные категории, такие как сыр, вино, печенье и молоко. Каждая категория поделена на подкатегории так, что продукты в одной подкатегории похожи и часто дублируют друг друга, но сами подкатегории существенно отличаются друг от друга. Например, категория «жидкое молоко» может включать подкатегории «цельное молоко», «обезжиренное молоко», «ароматизированное молоко» и т. д. Супермаркеты, как правило, достигают очень высокого уровня продаж продуктов, и затоваривание случается довольно редко, поэтому модель, которую мы рассмотрим, не учитывает возможность затоваривания. В то же время эта модель спроса разрабатывалась для анализа и оптимизации ассортимента, поэтому она явно учитывает потребительский выбор и хорошо подходит для решения задачи, связанной с ассортиментом, которую мы рассмотрим в последующих разделах.

Спрос на один продукт можно разбить на три составляющие, которые применяются к каждому потребителю, посещающему магазин.

- Во-первых, потребитель покупает или не покупает продукт из подкатегории. Обозначим вероятность приобретения (*purchase*) потребителем какого-либо продукта из подкатегории при посещении магазина (*visit*) как $\Pr(\textit{purchase} \mid \textit{visit})$.
- Во-вторых, потребитель выбирает продукт в подкатегории. Вероятность, что из всех альтернатив потребитель выберет продукт j , равна $\Pr(j \mid \textit{purchase})$.
- Наконец, потребитель решает, сколько единиц купить. Мы можем получить это число как математическое ожидание количества единиц продукта j , приобретенного потребителем, при условии, что этот продукт был выбран и куплен. Обозначим его как $\mathbb{E}[Q \mid j, \textit{purchase}]$.

Спрос на продукт j можно выразить через вероятности выбора и ожидаемое количество покупок:

$$D_j = N \times \Pr(\text{purchase} | \text{visit}) \times \Pr(j | \text{purchase}) \times \mathbb{E}[Q | j, \text{purchase}], \quad (6.50)$$

где N — количество потребителей, посещающих магазин в течение заданного периода времени (например, дня). Все члены в уравнении 6.50 можно определить по данным из магазинов о покупках. Спрос, как правило, зависит от даты (дня недели, праздников и т. д.) и магазина (размер, демография района и т. д.), поэтому введем подстрочные индексы t и h для обозначения даты и магазина соответственно и определим спрос как функцию этих параметров. Кроме того, характеристики магазина, такие как размер, местоположение и средний доход потребителя, можно включить в модель в качестве предиктивных переменных. Количество посетителей магазина можно смоделировать с помощью логарифмической линейной регрессии следующим образом:

$$\log(N_{ht}) = \alpha_1 + \alpha_2 T_t + \alpha_3 W_t + \sum_{i=1}^7 \alpha_{3+i} B_{ti} + \sum_{i=1}^{N_E} \alpha_{10+i} E_{ti}, \quad (6.51)$$

где T_t — температура воздуха на улице, W_t — индекс комфорта погоды (влажность, облачность и др.), B_{ti} и E_{ti} — принимают значение 0 или 1 в зависимости от дня недели и будних/праздничных дней соответственно, N_E — общее число праздничных дней и α — коэффициенты регрессии.

Факт покупки является переменной бинарного выбора (покупка/не покупка), поэтому можно использовать стандартный подход к моделированию — выразить вероятность покупки как сигмоидную функцию, аппроксимирующую бинарное решение, и определить ее экспоненциальный параметр по данным. Сигмоидную функцию можно задать как

$$\Pr(\text{purchase} | \text{visit}) = \frac{1}{1 + e^{-x}}, \quad (6.52)$$

что эквивалентно выражению

$$x = \log \left(\frac{\Pr(\text{purchase} | \text{visit})}{1 - \Pr(\text{purchase} | \text{visit})} \right). \quad (6.53)$$

Экспоненциальный параметр x оценивается для заданной даты t и магазина h с помощью регрессионной модели со следующей структурой:

$$x_{ht} = \beta_1 + \beta_2 T_t + \beta_3 W_t + \beta_4 \bar{A}_{ht} + \sum_{i=1}^7 \beta_{4+i} \beta_{ti} + \sum_{i=1}^{N_E} \beta_{11+i} E_{ti}, \quad (6.54)$$

где \bar{A}_{ht} — доля продуктов в подкатегории, которые продвигаются в данный момент, то есть отношение между количеством продуктов в подкатегории, которые продвигаются в данном магазине в данную дату, к общему количеству товаров в этой подкатегории. Так как в дальнейшем нам придется построить отдельную модель для каждого продукта, выразим доли продвигаемых продуктов в виде индикаторных переменных \bar{A}_{jht} , равных единице, если продукт продвигается, и ноль в противном случае:

$$\bar{A}_{ht} = \frac{1}{J} \sum_{j=1}^J A_{jht}, \quad (6.55)$$

где J — общее число продуктов в подкатегории.

Оценить вероятность покупки данного продукта в подкатегории несколько сложнее. Как мы видели выше, выбор потребителя можно смоделировать с помощью полиномиальной логит-модели (MNL), поэтому выразим вероятность выбора продукта из альтернатив следующим образом:

$$\Pr(j | purchase) = \frac{\exp(y_j)}{\sum_i \exp(y_i)}, \quad (6.56)$$

где i — индексы продуктов в подкатегории, а y_j — переменная параметра. По аналогии с вероятностью покупки можно построить регрессионную модель параметрической переменной y_j для данного магазина и даты:

$$y_{jht} = y_j + y_{N+1} (R_{jht} - \bar{R}_{ht}) + y_{N+2} (A_{jht} - \bar{A}_{ht}), \quad (6.57)$$

где коэффициенты регрессии γ_{N+1} и γ_{N+2} являются общими для всех продуктов, R_{jht} и \bar{R}_{ht} — цена товара и средняя цена в подкатегории соответственно, и A_{jht} и \bar{A}_{ht} — индикаторные переменные, определяющие факт продвижения, и средняя доля продвигаемых продуктов, как описывалось выше для регрессионной модели покупок.

Наконец, среднее количество проданных единиц можно смоделировать следующим образом:

$$\begin{aligned} \mathbb{E}[Q | j, \text{ purchase}] = & \lambda_j + \lambda_{N+1} A_{jht} + \lambda_{N+2} W_t + \\ & + \sum_{i=1}^{N_H} \lambda_{N+2+i} E_{ti}, \end{aligned} \quad (6.58)$$

где λ — коэффициенты регрессии, а другие переменные определяются, как было описано выше. Подставляя отдельные регрессионные модели в базовое уравнение 6.50, получаем полностью определенную модель прогнозирования спроса. Эту модель можно приспособить для конкретных бизнес-сценариев, добавив дополнительные переменные, например, определяющие маркетинговые события.

Конкурирующие продукты и их атрибуты играют важную роль в моделировании спроса, даже если ассортимент не является главной целью. Например, онлайн-магазин модной одежды Rue La La сообщил, что относительная цена конкурирующих стилей и их количество входят в тройку наиболее важных характеристик в модели прогнозирования спроса, используемой этим магазином [Ferreira et al., 2016].

6.6.2. Модель спроса для сезонных продаж

Вторая модель спроса, которую мы рассмотрим, была разработана для испанского ретейлера модной одежды Zara и головного бренда Inditex, крупнейшей в мире сети магазинов одежды [Caro and Gallien, 2012]. Модель ориентирована на оптимизацию распродаж и уделяет большое внимание измерению изменений спроса во времени.

Сезонные распродажи являются неотъемлемой частью бизнес-стратегии многих ретейлеров одежды. Сезоны продаж, которых обычно два в году (осень–зима и весна–лето), сопровождаются относительно короткими периодами распродаж, целью которых является продажа оставшихся запасов и освобождение места для новой коллекции на следующий сезон. Некоторые ретейлеры используют еще более короткие циклы продаж, чтобы обогнать конкурентов и получить больше доходов от клиентов, предлагая более разнообразный и гибкий ассортимент. Оптимизация цен в таких условиях требует создания модели спроса, которая учитывает сезонные эффекты и перебои, вызванные исчерпанием запасов и преднамеренными изменениями ассортимента.

Опишем модель спроса в два этапа, в соответствии с оригинальным отчетом [Caro and Gallien, 2012]. Первый шаг — подготовка имеющихся данных о спросе для регрессионного анализа устранением сезонных колебаний и учетом цензурирования спроса в связи с нехваткой запасов. Затем определим саму регрессионную модель.

6.6.2.1. Подготовка данных о спросе

Большинство предметов одежды имеют разные цвета и размеры, поэтому каждую единицу складского учета (Stock Keeping Unit, SKU) можно обозначить номером продукта r , а вариант размер/цвет как v . Предположив, что история продаж и данные о запасах доступны на уровне магазина для каждого дня, обозначим продажи единицы артикула $SKU(r, v)$ в магазине h за день d как $S(r, v, d, h)$ и уровень запасов на начало дня как $L(r, v, d, h)$. Определяем также функцию $F(r, v, d, h)$, равную единице, если данный артикул был доступен в магазине h в день d , и ноль в противном случае. Информация о доступности может присутствовать в данных явно или извлекаться из данных о продажах и запасах проверкой на равенство нулю уровня запасов или проданного количества данного продукта.

Во-первых, введем фактор сезонности, включающий внутри- и межнедельные колебания спроса, и определим следующие агрегаты данных о продажах:

- $S_w(d)$ — общий объем продаж за неделю, в которую попадает день d . Этот объем суммируется по всем продуктам, вариантам размеров/цветов и магазинам.
- \bar{S}_w — общий средний недельный объем продаж, рассчитанный на основе исторических данных.
- $\bar{S}_w(r)$ — средний недельный объем продаж продукта r , рассчитанный на основе исторических данных.
- $\bar{S}_w(r, v, h)$ — средний недельный объем продаж артикула $SKU(r, v)$ в магазине h , рассчитанный на основе исторических данных.
- $\bar{S}_{D(\text{weekday}(d))}$ — средний объем продаж в данный день недели. Определяется для семи дней — с понедельника по воскресенье.

Теперь фактор сезонности можно определить так:

$$\delta(d) = \frac{S_w(d)}{\bar{S}_w} \cdot \frac{\bar{S}_{D(\text{weekday}(d))}}{\sum_{i=1}^7 \bar{S}_D(i)}. \quad (6.59)$$

Первый и второй члены учитывают внутри- и межнедельные колебания спроса соответственно. Далее введем следующий фактор, учитывающий сезонность и наличие в продаже, чтобы нормализовать спрос для продукта r и недели w :

$$k(r, w) = \sum_{h,v} \frac{\bar{S}_w(r, v, h)}{\bar{S}_w(r)} \cdot \sum_{d \text{ in } w} \delta(d) \cdot F(r, v, d, h). \quad (6.60)$$

Уравнение выше определяет долю продаж в магазине h относительно всех магазинов, поэтому вклад переменной, описывающей наличие в продаже, правильно

взвешивается долей продаж магазина. Наконец, определим нормализованный спрос на продукт r и неделю w :

$$q(r, w) = \frac{1}{k(r, w)} \cdot \sum_{v, h, d \text{ in } w} S(r, v, d, h). \quad (6.61)$$

6.6.2.2. Определение модели

Наш следующий шаг — построение регрессионной модели, прогнозирующей нормализованный спрос $q(r, w)$. В Зага для этой цели использовали относительно небольшую логарифмически линейную модель со следующими характеристиками:

$$\begin{aligned} \log(q(r, w)) = & \alpha_{0,r} + \alpha_1 \log(Q_r) + \alpha_2 A_{r,w} + \\ & + \alpha_3 \log(q(r, w-1)) + \\ & + \alpha_{4,w} \log\left(\min\left\{1, \frac{1}{T} L(r, w)\right\}\right) + \\ & + \alpha_{5,w} \log\left(\frac{\text{Pr}, w}{\text{Pr}, 0}\right), \end{aligned} \quad (6.62)$$

где α — коэффициенты регрессии, и признаки определяются следующим образом:

- α_1 : Q_r — количество продукта r , приобретенное ретейлером. Хотя эта величина не связана напрямую со спросом, она косвенно связана с модными тенденциями и стилем, потому что ретейлер часто покупает большое количество популярных продуктов, тогда как нишевые продукты покупаются в меньшем количестве.
- α_2 : $A_{r,w}$ — количество дней с момента появления продукта r в магазинах. Спрос на модные изделия часто зависит от их новизны и имеет тенденцию к снижению с течением времени.
- α_3 : $q(r, w-1)$ — уровень спроса за предыдущий интервал времени. Эта переменная помогает выявить корреляцию спроса между соседними временными интервалами.
- $\alpha_{4,w}$: *эффект падения популярности ассортимента* — учитывает снижение спроса на данный продукт по мере уменьшения уровня запасов. В контексте торговли модной одеждой под этим часто подразумеваются непопулярные размеры и цвета, которые остаются после продажи самых популярных. Этот эффект можно учесть введением порога T для $L(r, w)$, который является агрегированным уровнем запасов для продукта r во всех вариантах и магазинах.

- $\alpha_{5,w}$: величина скидки, определяется как отношение текущей цены $p_{r,w}$ к обычной $p_{r,0}$. Этот член фактически является фактором ценовой чувствительности.

Как видите, модель в значительной степени ориентирована на изменчивость спроса с течением времени, потому что была создана для оптимизации сезонных распродаж. Модели спроса, о которых сообщают другие ретейлеры, торгующие модной одеждой, могут включать больше признаков, таких как популярность бренда, цвета и размера, относительные цены конкурирующих стилей и различные статистические данные о прошлых распродажах, которые могут изменить чувствительность к ценам [Ferreira et al., 2016].

6.6.3. Прогнозирование спроса с учетом истощения запасов

Модели спроса, описанные в предыдущих разделах, являются простыми регрессионными моделями, обученными для прогнозирования объемов продаж. На практике наблюдаемые объемы не обязательно соответствуют фактическому спросу из-за случаев отсутствия запасов. В такой ситуации наблюдаемый объем продаж будет ниже фактического спроса, то есть объема продаж, которого теоретически можно достичь при неограниченном предложении. Проблема запасов особенно актуальна для бизнес-моделей с сезонными или внезапными продажами, где запасы играют важную роль, и модель прогнозирования спроса, созданная на основе наблюдаемого объема продаж, почти наверняка окажется предвзятой и непригодной для оптимизации уровней запасов или цен. Поэтому необходимо разработать метод прогнозирования спроса, явно учитывающий случаи отсутствия продукта на складе и потерянные продажи. Эта задача активно изучалась, и существует ряд методов прогнозирования спроса с учетом истощения запасов [Anupindi et al., 1998; Musalem et al., 2010; Vulcano et al., 2012]. В этом разделе мы обсудим эвристический метод, разработанный в Rue La La — онлайн-магазине модной одежды — для учета случаев истощения запасов, возникающих во время внезапных распродаж, то есть в период действия скидок, чрезвычайно ограниченных по времени [Ferreira et al., 2016]. Преимущество этого метода заключается в простоте и способности работать с низким уровнем запасов (когда у ретейлера имеется лишь несколько экземпляров каждого артикула), то есть в условиях недостаточности информации для использования более продвинутых моделей.

Исходя из предположения, что ретейлер продает несколько продуктов, введем следующие обозначения:

- d_i — фактический спрос на продукт i ;
- c_i — уровень запасов продукта i ;

- q_i — фактический объем продаж, который можно выразить как

$$q_i = \min\{c_i, d_i\}. \quad (6.63)$$

Если продукт представлен несколькими вариантами размер/цвет, каждый рассматривается как отдельный продукт, и значения выше измеряются для каждого варианта. Теперь предположим, что ретейлер проводит ограниченные по времени распродажи отдельных продуктов. В начале распродажи уровень запасов равен c_i . Если продукт продан полностью до окончания распродажи, мы увидим, что $q_i = c_i$, но не увидим истинного спроса d_i . Если продукт не распродается, можно предположить, что мы наблюдаем истинный спрос, $q_i = d_i$. Основная задача, которую мы должны решить, заключается в том, как оценить ожидаемый объем продаж q_i с учетом уровня запасов c_i в качестве параметра. Можно выделить следующие случаи:

- Если продукт i уже имелся в продаже и не был продан, наблюдаемое проданное количество q_i можно использовать в качестве оценки спроса \hat{d}_i , а ожидаемое количество, которое будет продано с учетом уровня запасов c_i^{new} в новой распродаже, можно предсказать на основе этой оценки. Этот случай можно обобщить следующим образом:

$$\hat{d}_i = q_i \rightarrow \hat{q}_i = \min\{c_i^{new}, \hat{d}_i\}. \quad (6.64)$$

- Если продукт i уже имелся в продаже и был распродан, значит, истинного спроса мы не наблюдали. В этом случае мы должны определить *неограниченный спрос*, то есть оценить истинный спрос на основе проданного количества и исторических данных по другим продуктам. Эту оценку \hat{d}_i можно использовать для прогнозирования ожидаемого количества, которое будет продано в новой распродаже. Этот случай можно обобщить следующим образом:

$$q_i \text{ unconstraining } \hat{d}_i \rightarrow \hat{q}_i = \min\{c_i^{new}, \hat{d}_i\}. \quad (6.65)$$

- Если продукт новый и никогда прежде не продавался, спрос необходимо спрогнозировать с помощью регрессионной модели, использующей в качестве признаков свойства продукта и события, а оценку спроса — в качестве переменной отклика. Данную задачу можно решить с помощью методов, описанных в предыдущем разделе. Этот случай можно обобщить следующим образом:

$$\hat{d}_i = f(\text{product}, \text{event}) \rightarrow \hat{q}_i = \min\{c_i^{new}, \hat{d}_i\}, \quad (6.66)$$

где f — модель, прогнозирующая спрос. В случае с несколькими вариантами продукта можно создать модель для прогнозирования спроса на уровне продукта, а за-

тем на основе наблюдаемого распределения спроса для разных размеров и цветов вывести спрос на варианты.

Ключевой задачей в подходе, описанном выше, является определение величины неограниченного спроса. Один из возможных подходов к ее решению — использовать исторические данные по нераспроданным продуктам и на их основе оценить спрос на проданные продукты. Для начала можно создать кривые спроса для нераспроданных продуктов, чтобы каждая кривая описывала процент от общего объема продаж для данного продукта в зависимости от времени, измеренном в часах или днях, как показано на рис. 6.14.

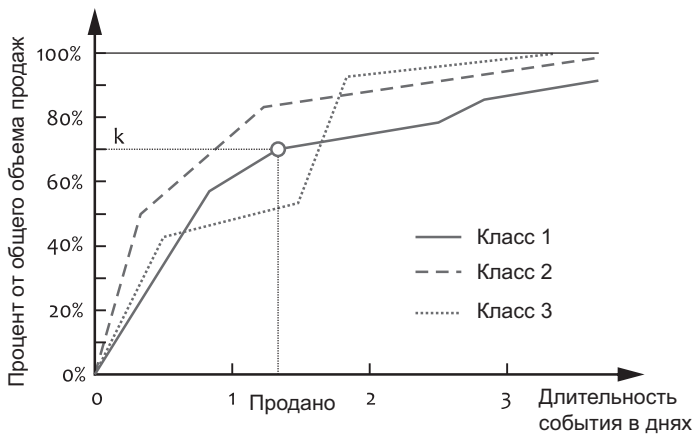


Рис. 6.14. Пример кривых спроса и определение доли неограниченного спроса [Ferreira et al., 2016]. В этом примере приводятся три класса событий, каждый из которых представлен типичной кривой спроса, и событие, для которого требуется оценить неограниченный спрос, попадает в первый класс

На практике число продуктов и соответствующих кривых спроса может быть очень большим, поэтому кривые можно разделить на несколько классов с помощью стандартных методов кластеризации, задав правила, которые отличают классы (кластеры) друг от друга, и определить типичную кривую спроса для каждого класса. Ключевым фактором, определяющим форму кривой спроса, может быть время начала события, например утро, ранний день, поздний день и т. д. Большинство событий, которые начинаются, например, утром, могут иметь похожие кривые спроса и, соответственно, принадлежать одному кластеру. Правило кластеризации может быть более сложным и включать несколько свойств. Определив правила, которые относят продукты и события к тем или иным классам, а также типичную кривую для каждого класса, мы сможем убрать ограничение спроса на продукты, которые

были проданы, на основе кривой. Сначала используем правила кластеризации, определяющие принадлежность продукта и события к некоторому классу. Затем используем кривую спроса для этого класса, чтобы сопоставить момент исчерпания запасов с долей общего объема продаж k , как показано на рис. 6.14. Определив эту пропорцию, истинный спрос на продукт i можно оценить, разделив наблюдаемое проданное количество на величину пропорции:

$$\hat{d}_i = \frac{q_i}{k} = \frac{c_i}{k}. \quad (6.67)$$

Эту оценку можно использовать для прогнозирования объема продаж, моделирования спроса и оптимизации уровня запасов.

6.7. Оптимизация цен

Модель спроса позволяет определить оптимальную цену, анализируя увеличение или уменьшение прибыли за счет изменения цены. Мы уже видели, что эта задача не особенно сложна для базовых ценовых структур и чрезмерно упрощенной среды, не учитывающей особые свойства спроса, предложения и операций, которые можно найти в реальных приложениях. Однако на практике часто приходится сталкиваться с многочисленными ограничениями и взаимозависимостями, которые требуют разработки более сложных и специализированных моделей оптимизации.

Большинство ограничений относятся к одной из следующих трех категорий: ограничения предложения, ограничения спроса и структурные ограничения. Ограничения предложения могут быть обусловлены ограниченностью ресурсов, таких как фиксированное число мест в театре, ограниченные возможности или высокая стоимость пополнения запасов и *скоропортящаяся* продукция, которые могут принимать разные формы — от ограниченного срока годности продуктовых товаров и сезонного характера коллекций одежды до авиабилетов, которые должны быть проданы до вылета самолета. Ограничения спроса часто связаны с несовершенной сегментацией потребителей, влиянием спроса на взаимозаменяемые продукты, изменениями спроса с течением времени и неопределенностью спроса в том смысле, что спрос нельзя точно предсказать. Структурные ограничения связаны с эксплуатационными и правовыми условиями, которые могут потребовать выбора неоптимальных, но практически осуществимых решений.

Мы продолжим этот раздел рассмотрением моделей оптимизации цен для различных стратегий ограничения и сегментации рынка. Эту группу методов можно рассматривать как статическую оптимизацию цен, однако процедуры оптимизации можно повторять регулярно для корректировки цен с течением времени или

определения сегментов на основе времени, таких как билеты в кино в выходные и будние дни. Затем мы увидим, что динамическое ценообразование, явно оптимизирующее изменения цен во времени, тесно связано с сегментацией рынка, и познакомимся с динамическими методами управления ценами.

6.7.1. Ценовая дифференциация

Цель ценовой дифференциации, которую в экономической литературе часто называют *ценовой дискриминацией*, заключается в определении оптимальных цен для отдельных сегментов потребителей или клиентов. Для оптимизации цен на уровне сегмента или клиента требуется создать модель спроса, которая принимает свойства клиента или сегмента в качестве параметров, или отдельную модель спроса для каждого сегмента. Эту задачу можно решить с помощью методов прогнозирования спроса, которые мы рассмотрели выше. Основную цель оптимизации цен можно определить следующим образом:

$$\max_p \sum_s (p_s - v_s) \cdot q_s(p_s), \quad (6.68)$$

где s — сегмент, p_s — цена для сегмента s , p — вектор цен для всех сегментов, q_s — функция спроса для сегмента s и v_s — переменные затраты, которые могут быть постоянными или изменяться в зависимости от сегмента. Эту задачу оптимизации можно разбить на сегменты, то есть базовую оптимизацию цены за единицу можно применить к каждому сегменту отдельно.

Часто число и структура сегментов ограничены эксплуатационными ограничениями, то есть программной системе может потребоваться оценить влияние объединения нескольких сегментов в группу и назначения этой группе единой цены. Это можно сделать, переписав уравнение 6.68 для N групп сегментов S_i как

$$\max_p \sum_{i=1}^N p_i \sum_{s \in S_i} (p_i - v_s) \cdot q_s(p_i), \quad (6.69)$$

и решить задачу оптимизации отдельно для каждой группы, чтобы найти N оптимальных цен.

ПРИМЕР 6.3

Проиллюстрируем модели оптимизации, описанные выше, на примере компании-ритейлера. Рассмотрим ритейлера, управляющего несколькими магазинами и продающего продукт, имеющий несколько размеров, напри-

мер обезболивающие таблетки в упаковках по 25 или 50 таблеток [Khan and Jain, 2005]. Ретейлер может предложить скидки в зависимости от размеров упаковки и установить цены отдельно для каждого магазина. Регрессионный анализ продаж показал, что спрос на анальгетики хорошо описывается следующей моделью:

$$q(p, s, h) = 2000 - 1400p - 8s - 10s \cdot h, \quad (6.70)$$

где p — цена одной таблетки, s — размер упаковки (количество таблеток) и h — коэффициент, учитывающий средний размер домохозяйства в районе расположения магазина: положительный, когда средний размер домохозяйства относительно велик, и отрицательный, когда средний размер домохозяйства мал. Спрос отрицательно коррелирует с ценой, что вполне ожидаемо. Он также отрицательно коррелирует с размером упаковки, то есть потребители предпочитают покупать более мелкие упаковки. Последний член положительно коррелирует с размером упаковки для больших домохозяйств и отрицательно — для малых, поэтому спрос на большие упаковки выше в районах с большими домохозяйствами, что также интуитивно понятно.

Наша задача — оптимизировать цены для случая с двумя магазинами в районах с разным значением коэффициента h , учитывающего средний размер домохозяйства, и упаковками двух размеров с разными оптовыми ценами v , как показано на рис. 6.15.

Первым рассмотрим сценарий дифференциации цен, оптимизирующий одновременно скидки за количество, на основе размеров упаковки, и цен на уровне магазина. Цель состоит в том, чтобы найти четыре разные цены p_{ij} , где i соответствует одному из двух размеров упаковки, а j — одному из двух магазинов. Задачу оптимизации можно сформулировать следующим образом:

$$\max_p \sum_{i=1,2} \sum_{j=1,2} (s_i p_{ij} - v_i) \cdot q(p_{ij}, s_i, h_j). \quad (6.71)$$

Эта задача оптимизации является separable и квадратичной по ценам, поскольку функция спроса является линейной. Решая задачу для значений из рис. 6.15, получим результаты, представленные в табл. 6.5. Как видите, решение оправдывает скидки за количество и предлагает установить более низкие цены на таблетки в больших упаковках в обоих магазинах. Оно также учитывает более высокий спрос на большие упаковки в районе с большими домохозяйствами, увеличивая цену в соответствующем магазине.

Второй сценарий предполагает невозможность варьировать цены на уровне магазина из-за эксплуатационных ограничений, поэтому мы можем установить только две цены — для маленьких и больших упаковок. Изменим задачу оптимизации в соответствии с общим подходом из уравнения 6.69:

$$\max_p \sum_{i=1,2} (s_i p_i - v_i) \sum_{j=1,2} q(p_i, s_i, h_j). \quad (6.72)$$



магазин 1

коэффициент среднего
размера домохозяйства:
 $h_1 = -0,1$



упаковка 1

размер: $s_1 = 25$ таблеток
оптовая цена: $v_1 = 6$ \$



магазин 2

коэффициент среднего
размера домохозяйства:
 $h_2 = 0,5$



упаковка 2

размер: $s_2 = 50$ таблеток
оптовая цена: $v_2 = 10$ \$

Рис. 6.15. Параметры в задаче оптимизации цен для двух магазинов и упаковок двух размеров

Таблица 6.5. Оптимальные цены для сценария с четырьмя сегментами

Размер упаковки (таблеток)	Размер домохозяйства	Цена за таблетку	Спрос (упаковок)
25	малый	\$0,67	607
25	большой	\$0,81	794
50	малый	\$0,49	410
50	большой	\$0,76	785
Общая выгода			\$45 863

Эта задача также является разделимой и квадратичной по ценам. Решая ее, находим, что общая прибыль уменьшается, по сравнению с первым сценарием, как можно видеть в табл. 6.6, где показаны результаты оптимизации.

Таблица 6.6. Оптимальные цены для сценария с двумя сегментами

Размер упаковки (таблеток)	Цена за таблетку	Спрос (упаковок)
25	\$0,74	1401
50	\$0,63	1195
Общая выгода		\$43 038

Этот анализ позволяет оценить различные стратегии сегментации и найти оптимальное решение при структурных ограничениях.

6.7.1.1. Дифференциация со смещением спроса

Одной из самых больших проблем ценовой дифференциации является несовершенство барьеров между ценовыми сегментами, что позволяет клиентам переходить из одного сегмента в другой в зависимости от разницы цен. Звучит так, будто такое смещение спроса обязательно вредно для продавца, но на самом деле это может влиять на прибыль как положительно, так и отрицательно. С одной стороны, клиенты с высокой готовностью платить могут найти способ купить продукт по относительно низкой цене, тем самым снижая прибыль высокодоходных сегментов. С другой стороны, увеличение сегмента, в который перешел клиент, может компенсировать потери.

Эффект смещения особенно важен в случаях ограниченного предложения, потому что помогает добиться более равномерного распределения спроса и сократить запасы. Например, количество мест в театре фиксировано, но спрос может меняться значительно, обычно достигая максимума в выходные дни и минимума в будни. Театр может потерять потенциальную выручку, если пиковый спрос в выходные превысит пропускную способность. Можно ожидать, что установление переменных цен на билеты в разные дни недели улучшит прибыль, потому что более высокий спрос в выходные позволяет театру взимать более высокую плату и получать лучшую маржу, как было показано в предыдущем разделе. Однако высокие цены в выходные могут заставить некоторых клиентов покупать более дешевые билеты в будни, что сдвигает спрос на дни с меньшей заполненностью зала и улучшает доходы.

Создать модель смещения спроса порой очень непросто. Сложность отчасти связана с разработкой и обучением модели спроса, которая одновременно учитывает цены всех связанных сегментов. Часто невозможно измерить кросс-ценовую эластичность (влияние цены в одном сегменте на спрос в другом) для всех возможных пар сегментов, поэтому приходится использовать более грубые приближения, такие

как соотношение между ценой в данном сегменте и средними ценами в других сегментах. Еще одна сложность построения моделей смещения спроса состоит в том, что межсегментные зависимости делают задачу оптимизации неразделимой и резко увеличивают ее вычислительную сложность. Если допустить, что цены выбираются из дискретного множества размера m , а количество сегментов равно n , нам может потребоваться оценить m^n комбинаций цен, если модель спроса не следует какой-то конкретной функциональной форме, проявляющей такие свойства, как линейность или выпуклость.

Один из возможных подходов к оптимизации цен со смещением спроса заключается в предположении, что смещение спроса пропорционально разнице цен между сегментами. То есть если цена в сегменте i выше, чем в сегменте j , то спрос в сегменте i уменьшается на $K(p_i - p_j)$, а в сегменте j — увеличивается на эту величину. Параметр K определяет величину спроса, передаваемого между двумя сегментами на каждый доллар разницы в цене. Основную задачу оптимизации в этом случае можно переписать, как показано ниже, чтобы скорректировать спрос в каждом сегменте на сумму спроса, сдвинутого из других сегментов:

$$\max_p \sum_i (p_i - v_i) \cdot \left[q_i(p_i) + K \sum_j (p_j - p_i) \right], \quad (6.73)$$

где i и j перебирают все сегменты. Обратите внимание, что эта модель смещения спроса не меняет общей величины спроса, то есть сумма всех смещений равна нулю. Но это не означает, что общее проданное количество остается постоянным для любого K , потому что смещение спроса приводит к изменению оптимальных цен, что, в свою очередь, изменяет значения функций спроса.

ПРИМЕР 6.4

Проиллюстрируем влияние смещения спроса, продолжив пример с обезболивающими таблетками. Вставив члены, определяющие смещение спроса, в уравнение 6.71, которое описывает сценарий с четырьмя ценовыми сегментами, получаем:

$$\max_p \sum_{i=1,2} \sum_{j=1,2} (s_i p_{ij} - v_i) \cdot \left[q(p_{ij}, s_i, h_j) + \Delta(p_{ij}) \right], \quad (6.74)$$

где смещение спроса является суммой попарных разниц цен с другими сегментами, или

$$\Delta(p_{ij}) = K \sum_{x=1,2} \sum_{y=1,2} (p_{xy} - p_{ij}). \quad (6.75)$$

Решив эту задачу оптимизации, получим цены для четырех сегментов, представленные в табл. 6.7. Если сравнить результаты табл. 6.5 и 6.7, можно обнаружить, что смещение спроса повысило чувствительность к ценам, поэтому спрос на маленькие упаковки снизился, а на большие — увеличился из-за относительно низкой цены за таблетку. Такое изменение спроса можно считать положительным для ретейлера, потому что увеличивает общую прибыль.

Таблица 6.7. Оптимальные цены для сценария с четырьмя сегментами и смещением спроса. Параметр смещения $K = 400$

Размер упаковки (таблеток)	Размер домохозяйства	Цена за таблетку	Спрос (упаковок)
25	малый	\$0,75	420
25	большой	\$0,81	608
50	малый	\$0,56	535
50	большой	\$0,69	910
Общая выгода			\$46 170

6.7.1.2. Дифференциация при ограниченном предложении

Модели оптимизации цен, которые мы рассматривали до сих пор, сосредоточены на определении цен, позволяющих достичь максимальной прибыли, допускаемой кривой спроса. Такой взгляд на оптимизацию цен предполагает идеальное пополнение запасов, когда продавец всегда может предложить требуемое количество продукта по цене, оптимальной с точки зрения прибыли. Это вполне приемлемое допущение, например, для супермаркета, когда можно построить цепочку поставок, которая почти идеально пополняет запасы, и ситуация исчерпания запасов возникает очень редко. Однако, как уже говорилось, это не относится ко многим другим сферам, где могут иметь место различные ограничения предложения. В этом разделе мы рассмотрим относительно простой случай, когда каждый сегмент рынка имеет фиксированную емкость продукта и необходимо найти оптимальные глобальные цены или цены на уровне сегмента.

Для начала посмотрим, как можно оценить продукт для одного маркетингового сегмента, если доступное его количество фиксировано. Это стандартная задача оптимизации цены за единицу с добавлением ограничения количества:

$$\begin{aligned}
 & \max_{p, x} \quad x(p - V) \\
 & \text{при условии} \quad x \leq q(p) \\
 & \quad \quad \quad x \leq C \\
 & \quad \quad \quad x \geq 0,
 \end{aligned}
 \tag{6.76}$$

где C — доступное количество (емкость), то есть если спрос превышает C , возникает событие исчерпания запасов. Как обычно, обозначим цену, спрос и переменные затраты как p , $q(p)$ и V соответственно. Переменная x соответствует фактически проданному количеству.

Эта задача достаточно тривиальна, потому что спрос является монотонно убывающей функцией цены. Прежде всего, можно найти оптимальную цену без учета ограничения и рассчитать спрос, соответствующий этой неограниченной оптимальной цене. Затем можно сравнить этот спрос с доступным количеством и установить цену, исходя из максимума этих двух значений. Если спрос, соответствующий неограниченной оптимальной цене, ниже доступного количества, тогда эта цена является решением, потому что уровень запасов является менее строгим ограничением, чем спрос. Иначе берем цену, соответствующую максимальному доступному количеству. Эта цена, называемая *ценой дефицита*, должна быть выше неограниченной оптимальной цены, чтобы замедлить продажи и избежать исчерпания запасов (дефицита). Последний случай показан на рис. 6.16.

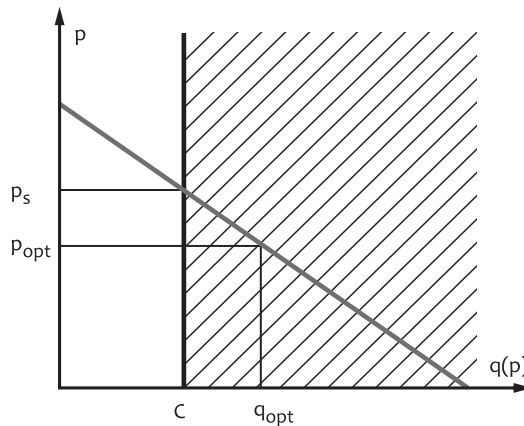


Рис. 6.16. Оптимизация цены за единицу продукта с ограниченным предложением

Описанную идею можно объединить с приемом дифференциации цен со смещением спроса, рассмотренными выше. Однако такой подход применим, только

если ограничения предложения установлены для каждого сегмента отдельно. Глобальные ограничения емкости, когда все ценовые сегменты имеют одинаковый запас, требуют применения более продвинутых методов оптимизации, которые мы рассмотрим далее в этой главе.

ПРИМЕР 6.5

Возьмем в качестве примера театр и попробуем оптимизировать цены на билеты. Предположим, что театр дает спектакли ежедневно, вместимость (C) зрительного зала составляет 1200 мест, а спрос меняется в течение недели в соответствии со следующими формулами:

$$q(p, t) = \begin{cases} 1800 - 50p, & \text{понедельник} \\ 1350 - 50p, & \text{вторник} \\ 1200 - 50p, & \text{среда} \\ 1350 - 50p, & \text{четверг} \\ 1800 - 50p, & \text{пятница} \\ 2250 - 50p, & \text{суббота} \\ 3600 - 50p, & \text{воскресенье} \end{cases}, \quad (6.77)$$

где p — цена билета, а t — день недели. Предположим также, что переменные затраты за место незначительны, потому что затраты на техническое обслуживание здания, на подготовку спектакля и другие расходы в значительной степени постоянны.

Руководство театра может решить установить фиксированную цену на все дни, что приводит нас к следующей задаче ограниченной оптимизации:

$$\begin{aligned} \max_{p, x} \quad & p \sum_t x_t \\ \text{при условии} \quad & x_t \leq q(p, t) \\ & x_t \leq C \\ & p \geq 0, \end{aligned} \quad (6.78)$$

где t перебирает семь дней недели. Эта задача решается достаточно просто, потому что можно предположить, что цена принадлежит относительно небольшому дискретному множеству, и оценить каждое потенциальное

решение. Решение этой задачи дает оптимальную цену 19,80 долларов, что соответствует доходу 98 010 долларов.

Полученный результат можно сравнить с переменным ценообразованием, когда каждый день рассматривается и оптимизируется как отдельный сегмент. Перепишем задачу оптимизации для семи разных цен:

$$\begin{aligned} \max_{p, x} \quad & \sum_t p_t x_t \\ \text{при условии} \quad & x_t \leq q(p_t, t) \\ & x_t \leq C \\ & p_t \geq 0. \end{aligned} \quad (6.79)$$

Эта задача делится на сегменты, поэтому каждый день можно оптимизировать в соответствии с уравнением 6.76. Оптимальные цены, а также количество проданных мест x_t , показаны на рис. 6.17, наглядно демонстрирующем увеличившуюся загрузку зала. Дифференциация цен в этом примере оказалась чрезвычайно эффективной и увеличила общий доход до 147 825 долларов.

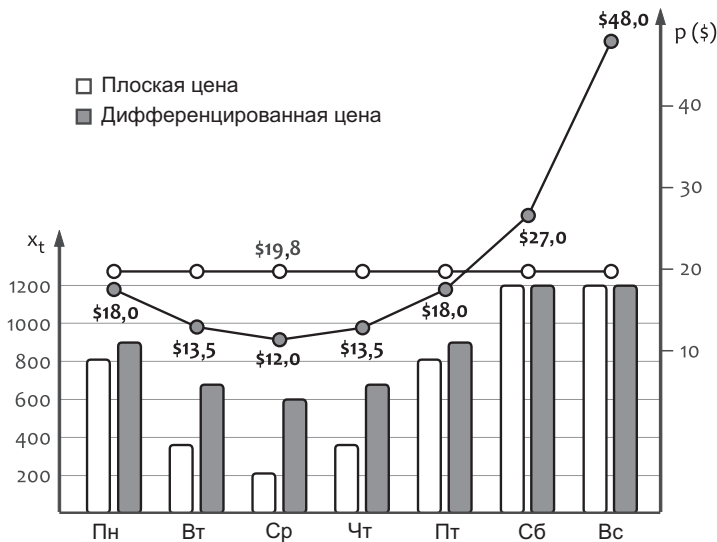


Рис. 6.17. Пример оптимизации цен на билеты в театр. Вертикальные столбики представляют число проданных мест, а точки — цены на билеты в соответствующие дни

6.7.2. Динамическое ценообразование

Под динамическим ценообразованием подразумевается группа стратегий и методов, оптимизирующих прибыль за счет изменения цены с течением времени. Как будет показано ниже, эти методы структурно аналогичны другим методам дифференциации цен, но предлагают дополнительные преимущества для алгоритмического подхода. В частности, динамическое ценообразование в большей степени ориентировано на поэтапные, постоянно корректируемые ценовые решения, которые делают решение более гибким, чем глобальная оптимизация цен. Кроме того, динамическое ценообразование, как правило, учитывает строгие ограничения емкости и времени, которые вводят дополнительные сложности, что, в свою очередь, увеличивает привлекательность автоматизированной оптимизации.

Интуитивно понятно, что динамическое ценообразование обеспечивает дополнительную степень свободы, по сравнению со статическими ценами, поэтому мы должны глубже изучить эту идею, чтобы понять, почему изменение цены с течением времени может приносить дополнительную выгоду.

Прежде всего отметим, что динамическое ценообразование можно использовать как метод сегментации рынка. Рассмотрим продавца, который изначально предлагает продукт по базовой цене. Продукт будет приобретаться клиентами, чья готовность платить выше базовой цены, и не будет приобретаться другими клиентами. Предположив, что все клиенты приняли решение в отношении возможности или невозможности покупки по базовой цене, продавец может снизить цену. Это позволит ему получить дополнительный доход от покупателей с готовностью платить между базовой ценой и ценой со скидкой. То есть динамическое ценообразование способно создавать ценовые сегменты только на основе времени и отклонений в готовности платить, без возведения дополнительных барьеров между сегментами. Эту идею мы используем в следующем разделе для построения количественной модели, которая позволяет провести оптимальную траекторию цены.

Неоднородность готовности платить может быть достаточным условием для динамического ценообразования, но в значительном числе бизнес-моделей спрос демонстрирует еще большую изменчивость. В таких ситуациях динамическое ценообразование играет роль регулятора, корректирующего цены в соответствии с изменяющимися условиями спроса.

ИЗМЕНЧИВЫЙ СПРОС. Спрос на продукт или услугу может меняться с течением времени, часто в соответствии с сезонными закономерностями. К отраслям, сталкивающимся с изменчивым спросом, можно отнести ретейлеров, торгующих одеждой, развлекательные предприятия и отели.

ИЗМЕНЧИВАЯ ЦЕННОСТЬ. Изменения спроса часто связаны с объективными или субъективными изменениями ценности продуктов. Модные продукты, электронные устройства и автомобили теряют свою ценность с появлением на рынке новых моделей. Скоропортящиеся продукты питания теряют ценность по мере приближения к границе срока годности, тогда как авиабилеты, как правило, более ценны для покупателей, принимающих решение в последнюю минуту.

НЕОПРЕДЕЛЕННОСТЬ СПРОСА. Динамическое изменение цен может помочь найти оптимальную цену методом проб и ошибок, когда объем спроса заранее не известен продавцу [Pashigian, 1987]. Например, ретейлер, торгующий одеждой и покупающий ее заранее, перед следующим сезоном, может не суметь предсказать популярность нового предмета одежды. Однако он может попробовать различные скидки, чтобы найти цену, максимизирующую прибыль и соответствующую ограничениям запасов.

С точки зрения оптимизации, динамическое ценообразование не всегда требует применения специализированных методов. Предположив, что время дискретно, когда каждый временной интервал можно рассматривать как сегмент и нет ограничений, делающих интервалы взаимозависимыми, мы можем оптимизировать интервалы по отдельности, используя методы, аналогичные случаям дифференциации цен (вспомните пример с театром). Зависимости между временными интервалами могут потребовать создания специализированных моделей оптимизации, что часто имеет место в динамическом ценообразовании. Эти зависимости обычно связаны с ограничениями предложения, поскольку все временные интервалы обычно обслуживаются из одного пула продуктов.

Выше уже отмечалось, что динамическое ценообразование можно рассматривать как регулятор спроса. Это рассуждение можно продолжить и сказать, что динамическое ценообразование регулирует не только спрос, но и предложение. Два главных атрибута, ограничивающих предложение, — фиксированная емкость и склонность к порче.

ФИКСИРОВАННАЯ ЕМКОСТЬ. Глобальное ограничение емкости означает, что продавец продает фиксированный запас продукта, который нельзя пополнить. Хорошим примером является розничная продажа одежды, когда ретейлер покупает фиксированные объемы запасов заранее. Сюда же можно отнести случай, когда пополнение возможно, но в ограниченных объемах.

СКЛОННОСТЬ К ПОРЧЕ. Под склонностью к порче подразумевается необходимость продажи запаса продукта в течение ограниченного времени. Непроданные запасы теряют свою стоимость или могут быть проданы лишь по относительно небольшой цене. Примерами продуктов, склонных к порче, могут служить ресурсы

обслуживания, такие как гостиничные номера и авиабилеты, сезонные коллекции одежды и потребительские упакованные товары.

Наличие этих двух ограничений в бизнесе, как правило, является хорошей предпосылкой для динамического ценообразования. Эти два свойства одинаково важны, потому что наша цель — соотнести величину спроса с величиной предложения, что по сути является отношением емкости к периоду времени продажи, определяемому из срока годности.

6.7.2.1. Уценки и распродажи

Динамическую оптимизацию цен с фиксированной емкостью скоропортящихся запасов можно выразить в виде следующей математической задачи:

$$\begin{aligned} \max_{p, x} \quad & \sum_{t=1}^T p_t x_t \\ \text{при условии} \quad & \sum_{t=1}^T x_t \leq C \\ & x_t \leq q(p_t, t), \quad \text{для } t = 1, \dots, T \\ & p_t \geq 0 \quad \text{для } t = 1, \dots, T. \end{aligned} \tag{6.80}$$

В этой формулировке предполагается, что объем C должен быть распродан в течение периода времени, состоящего из T дискретных временных интервалов. Наша цель — максимизировать доход, устанавливая оптимальные цены для каждого из T временных интервалов. Также предполагается, что нереализованные запасы имеют нулевую стоимость после точки T , а переменные затраты незначительны, впрочем, в уравнение легко можно включить и стоимость утилизации запасов, и переменные затраты.

Задача 6.80 моделирует несколько важных бизнес-сценариев. В сфере розничной торговли этой модели соответствуют уценки и сезонные распродажи, потому что период продажи обычно фиксирован, как и объем, подлежащий продаже. С концептуально схожими задачами оптимизации сталкиваются разнообразные поставщики услуг, включая авиакомпании, железные дороги, гостиницы, театры, стадионы и грузовые компании, при продаже фиксированного числа мест или номеров в течение ограниченного периода времени, определяемого расписанием движения, датой регистрации или временем события. Однако сфера услуг часто сталкивается со многими дополнительными ограничениями и использует различные методы оптимизации доходов, основанные на распределении ресурсов, поэтому в этом разделе мы сосредоточимся на уценках и распродажах. Методы распределения ресурсов мы рассмотрим далее в этой главе.

Первый вывод, который можно сделать из уравнения 6.80, — менять цены с течением времени можно, только если меняется уровень спроса. Если спрос находится на постоянном уровне, тогда все временные интервалы идентичны и можно применить стандартный прием оптимизации цены за единицу, а затем выбрать максимум этой неограниченной оптимальной цены и цены запаса, определяемый ограничением емкости C , согласно логике, описанной в разделе 6.7.1.2.

Изменчивость спроса с течением времени может объясняться различными факторами, такими как сезонность или изменение ценности продукта. Изменения спроса с течением времени могут принимать различные формы, включая тенденции к увеличению и уменьшению, как и цена. Однако мы можем показать, что уценки часто выполняются по определенной схеме и с учетом ограниченности круга клиентов [Talluri and Van Ryzin, 2004]. Предположим, что некто продает долговечный продукт конечному числу клиентов, и каждый клиент покупает продукт в течение периода продажи только один раз. Если продавец установит определенную цену p_t в интервале времени t , то все покупатели с готовностью заплатить цену выше или равную p_t купят продукт и станут неактивными до конца периода продажи. Однако продавец может снизить цену, чтобы привлечь клиентов с меньшей готовностью платить. Такая стратегия означает, что в начале продаж цены должны устанавливаться как можно ближе к оценке клиента, а затем монотонно уменьшаться.

Для построения количественной модели предположим также, что готовность клиентов платить равномерно распределена между 0 и некоторой максимальной ценой P :

$$w(p) = \text{unif}(0, P) = \begin{cases} 1/P, & 0 \leq p \leq P \\ 0 & \text{в противном случае.} \end{cases} \quad (6.81)$$

Напомним, что равномерное распределение готовности платить подразумевает линейную кривую спроса, которую в данном контексте можно интерпретировать как количество покупателей, покупающих товар по заданной цене и становящихся неактивными до конца продажи. Следовательно, процесс уценки можно изобразить как скольжение вниз по кривой спроса, как показано на рис. 6.18. Обратите внимание, что оптимизация уценок в данной интерпретации практически идентична задаче сегментации рынка, которую мы изучали выше, поэтому уравнения ниже структурно схожи с уравнениями сегментации рынка, но имеют иной смысл.

Итак, количество проданного продукта за период t выражается как

$$\begin{aligned} Q_t &= Q_{\max} \left[\left(1 - \frac{p_t}{P} \right) - \left(1 - \frac{p_{t-1}}{P} \right) \right] = \\ &= \frac{Q_{\max}}{P} (p_{t-1} - p_t). \end{aligned} \quad (6.82)$$

Следовательно, общий доход от реализации составит

$$G = \sum_{t=1}^T p_t \frac{Q_{\max}}{P} (p_{t-1} - p_t). \quad (6.83)$$

Возьмем частные производные, чтобы найти цены, максимизирующие доход:

$$\frac{\partial G}{\partial p_t} = \frac{Q_{\max}}{P} (p_{t-1} - 2p_t + p_{t+1}). \quad (6.84)$$

Приравняв эти производные к нулю, получим

$$p_t = \frac{p_{t-1} + p_{t+1}}{2}. \quad (6.85)$$

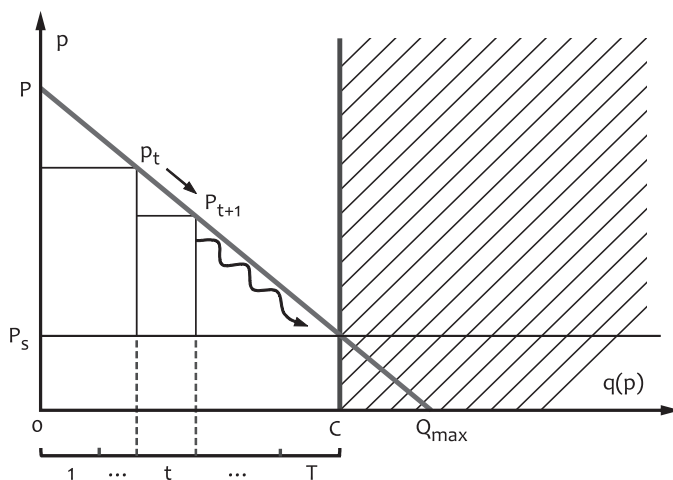


Рис. 6.18. Модель оптимизации цен для равномерно распределенной готовности платить

Начальные условия для этого рекуррентного уравнения должны устанавливаться в соответствии с ограничениями емкости, чтобы цены распределялись в диапазоне между максимальной ценой P и ценой запаса P_s :

$$\begin{aligned} p_0 &= P \\ p_{T+1} &= P_s. \end{aligned} \quad (6.86)$$

Цены, отвечающие отношению 6.85 и условиям 6.86, определяются как

$$p_t^{opt} = P_s + (P - P_s) \left(1 - \frac{t}{T+1} \right). \quad (6.87)$$

Этот результат предполагает, что цены должны равномерно распределяться между исходной ценой и ценой запаса. Это очень интересное концептуальное представление, но модель ограниченного круга клиентов слишком грубая для практических ценовых решений, поэтому вернемся к основной задаче оптимизации 6.80 и решим ее непосредственно для произвольных функций спроса. Такой подход даст большую гибкость и позволит учитывать в моделях прогнозирования спроса различные закономерности, наблюдавшиеся в прошлом, включая эффекты, связанные с ограниченностью круга клиентов.

6.7.2.2. Оптимизация уценки

Задача оптимизации 6.80 играет очень важную роль в алгоритмическом ценообразовании, поэтому посвятим этот раздел изучению ее эффективного решения на численном примере.

Предположим, что множество допустимых цен дискретно, это верно в большинстве практических приложений, где цены выражаются целыми значениями центов или долларов. Если множество допустимых цен имеет K уровней, а количество раундов уценки равно T , тогда оптимальные цены для каждого раунда можно найти, оценив K^T возможных ценовых комбинаций. Такой подход осуществим для случаев с небольшими величинами K и T . Например, многие ретейлеры, практикующие уценку, часто устанавливают цены, заканчивающиеся на 4.90 или 9.90, например 34.90 или 59.90 долларов. Однако проблема становится неразрешимой, даже если есть всего двадцать ценовых уровней и требуется оптимизировать ежедневные цены на горизонте двух недель, потому что это потребует оценить 20^{14} ценовых комбинаций.

Одним из возможных решений является аппроксимация исходной нелинейной задачи оптимизации с применением ослабленной задачи линейного программирования [Talluri and Van Ryzin, 2004]. Обозначим множество допустимых ценовых уровней как $\{P_1, \dots, P_K\}$ и введем весовые переменные z , управляющие ценовыми уровнями, выбранными для каждого временного интервала так, что

$$p_t = \sum_{i=1}^K z_{it} P_i. \quad (6.88)$$

Теперь задачу оптимизации можно переформулировать как поиск $K \times T$ переменных z , максимизирующих доход, таких, что для любого временного интервала t

переменная z_{it} равна единице, а другие переменные z равны нулю. Давайте ослабим ограничение «ноль или единица» для переменных z и просто потребуем, чтобы их сумма была равна единице для каждого временного интервала:

$$\sum_{i=1}^K z_{it} = 1, \quad \text{для } t = 1, \dots, T$$

$$z_{it} \geq 0.$$
(6.89)

Это фактически означает допустимость дробных цен, то есть две или более цены из дискретного множества могут иметь ненулевые веса в одном временном интервале. Мы предполагаем, что это приемлемо, а значит, раунд уценки можно разбить на более мелкие интервалы. Например, в случае двух ценовых уровней с ненулевыми весами 0,2 и 0,8 первый должен распространяться на одну пятую продолжительности раунда, а второй — на оставшиеся четыре пятых. Такое ослабление позволяет переписать задачу динамической оптимизации цены следующим образом:

$$\max_z \quad \sum_{t=1}^T \sum_{i=1}^K z_{it} \cdot P_i \cdot q(P_i, t)$$

при условии

$$\sum_{t=1}^T \sum_{i=1}^K z_{it} \cdot q(P_i, t) \leq C$$

$$\sum_{i=1}^K z_{it} = 1, \quad \text{для } t = 1, \dots, T$$

$$z_{it} \geq 0.$$
(6.90)

Задача 6.90 является задачей линейного программирования: если поместить все переменные z в плоский вектор с $K \times T$ элементами и вычислить векторы значений доходов и спроса для соответствующих пар времени t и уровня цены i , тогда целевую функцию и ограничение емкости можно выразить в виде векторных произведений. Эта формулировка позволяет использовать стандартное ПО оптимизации методом линейного программирования.

ПРИМЕР 6.6

Рассмотрим пример ретейлера, планирующего четырехнедельную кампанию по продаже одного продукта. Множество допустимых цен включает пять уровней: 89, 79, 69, 59 и 49 долларов. Функции спроса для каждой недели оцениваются следующим образом:

$$q(p, t) = \begin{cases} 1800 - 10p, & \text{неделя 1} \\ 1300 - 15p, & \text{неделя 2} \\ 1200 - 15p, & \text{неделя 3} \\ 1100 - 18p, & \text{неделя 4.} \end{cases}$$
(6.91)

Подставим эти параметры в задачу 6.90 и решим ее для разных значений емкости C . Полученные решения приводятся в таблицах ниже. Каждое решение представляет ценовой график с 20 значениями z : каждый столбец соответствует одной из четырех недель, а каждая строка — одному из пяти уровней, причем самая верхняя строка соответствует самой высокой цене, а самая нижняя строка — самой низкой. Например, оптимальная цена составляет 89 долларов в первую, третью и четвертую недели, когда емкость составляет 700 единиц. Во вторую неделю для той же емкости цена получила дробное значение, что означает смесь цен 89 и 79 долларов в пропорции 22 % и 78 % соответственно.

$C = 700, G = \$61,400$				$C = 1000, G = \$81,507$			
1,00	0,22	1,00	1,00	1,00	0,00	0,00	1,00
0,00	0,78	0,00	0,00	0,00	0,27	0,00	0,00
0,00	0,00	0,00	0,00	0,00	0,73	1,00	0,00
0,00	0,00	0,00	0,00	0,00	0,00	0,00	0,00
0,00	0,00	0,00	0,00	0,00	0,00	0,00	0,00
$C = 1300, G = \$96,778$				$C = 1600, G = \$109,218$			
1,00	0,00	0,00	0,00	1,00	0,00	0,00	0,00
0,00	0,00	0,00	0,00	0,00	0,00	0,00	0,00
0,00	1,00	0,00	0,00	0,00	0,00	0,00	0,00
0,00	0,00	1,00	0,60	0,00	1,00	0,72	0,00
0,00	0,00	0,00	0,40	0,00	0,00	0,28	1,00
$C = 1600, G = \$116,198$				$C = 2200, G = \$117,242$			
0,58	0,00	0,00	0,00	0,00	0,00	0,00	0,00
0,42	0,00	0,00	0,00	1,00	0,00	0,00	0,00
0,00	0,00	0,00	0,00	0,00	0,00	0,00	0,00
0,00	0,00	0,00	0,00	0,00	0,00	0,00	0,00
0,00	1,00	1,00	1,00	0,00	1,00	1,00	1,00

Как видите, цены в целом со временем снижаются, что отражает тенденцию к снижению спроса. Другая тенденция заключается в снижении жесткости ограничения емкости и замедлении снижения цены, что также вполне ожидаемо.

6.7.2.3. Оптимизация цен на конкурирующие продукты

Одной из основных проблем оптимизации цен являются зависимости между продуктами. Во многих сферах бизнеса, особенно в розничной торговле, клиенты постоянно выбирают между конкурирующими или взаимозаменяемыми продуктами, используя цену одного продукта в качестве ориентира для другого продукта, поэтому спрос на данный продукт обычно зависит не только от его цены, но также от цен на конкурирующие продукты. В этом случае цены не могут оптимизироваться для каждого продукта в отдельности — цены на все конкурирующие продукты должны оптимизироваться совместно. В некоторых случаях количество конкурирующих продуктов может достигать нескольких сотен, поэтому задача оптимизации может стать вычислительно неразрешимой. В этом разделе мы углубимся в эту проблему и обсудим структуру, разработанную Rue La La, онлайн-магазином модной одежды, которая может значительно уменьшить объем вычислений [Ferreira et al., 2016].

Предположим, что ретейлер продает n разных продуктов. Обозначим множество продуктов как N , то есть $|N| = n$, и множество возможных уровней цен как P , то есть $|P| = k$. На практике множество цен часто относительно невелико, потому что обычно принято использовать цены с определенной окантовкой, например 9.95 и 14.95. Предположим также, что ретейлер имеет модель прогнозирования спроса, которая оценивает спрос на данный продукт по его цене и ценам на конкурирующие продукты. В случае неограниченных запасов модель может предсказать истинный (неограниченный) спрос. Если ретейлер имеет фиксированное количество единиц продукта на складе, модель может предсказать ожидаемый объем продаж, которое является минимумом истинного спроса и доступного количества, как говорилось выше, в разделе 6.6.3.

Целью оптимизации является назначение такой цены каждому продукту, чтобы общий доход оказался максимальным. Данная постановка задачи применима как в случае статического, так и динамического ценообразования. В случае динамического ценообразования, например, нескольких параллельных распродаж, эту задачу оптимизации необходимо решать многократно. Например, ретейлер может инициировать n распродаж для всех n продуктов одновременно с определенным количеством каждого продукта на складе. Начальные цены выбираются путем решения задачи оптимизации для начальных настроек. На следующий день цены выбираются повторным решением задачи с последними уровнями запасов, и так далее. Простейшее решение — попробовать все k^n возможных цен и оценить спрос и доходы для каждой. Формально эту задачу можно определить так:

$$\begin{aligned} \max_p \quad & \sum_{i \in N} p_i \cdot q_i(p) \\ \text{при условии} \quad & p_i \in P, \text{ для } i = 1, \dots, n, \end{aligned} \tag{6.92}$$

где \mathbf{p} — n -мерный вектор цен на продукт, p_i — цена на продукт i , а q_i — модель спроса на продукт i , которая использует в качестве входных данных все цены, включая цену на продукт p_i и цены на все конкурирующие продукты. Количество продуктов n и количество возможных цен k могут быть относительно большими, поэтому такой подход непригоден для многих практических приложений. Обойти эту проблему можно, если реализовать прогнозирование спроса не как функцию всех цен на отдельные продукты, а как функцию некоторого агрегата с меньшим числом возможных состояний, чем вектор \mathbf{p} . Например, в качестве агрегата можно использовать сумму цен на конкурирующие товары:

$$Q = \sum_{i=1}^n p_i. \quad (6.93)$$

Легко заметить, что общее число m возможных значений Q составляет

$$m = n(k-1) + 1. \quad (6.94)$$

Например, если количество товаров равно 10, а множество возможных цен включает 10 уровней $P = \{\$1, \dots, \$10\}$, то значение Q будет находиться в диапазоне от 10 до 100 долларов с шагом 1 доллар. Предположение, что спрос зависит главным образом от суммы, а не от отдельных цен, на практике может оказаться ошибочным, но есть доказательства, что этот подход работает, по крайней мере, в некоторых приложениях [Ferreira et al., 2016]. Это предположение существенно сокращает пространство поиска от k^n комбинаций цен до $O(nk)$ возможных значений Q . Оптимальную цену в этом случае можно найти путем решения m задач оптимизации для каждого возможного значения Q и выбора наилучшего результата. Чтобы определить задачу оптимизации, обозначим j -й элемент множества ценовых уровней P как $p^{(j)}$ и введем бинарные переменные $z_{ij} \in \{0, 1\}$, такие, что z_{ij} равна единице, если продукту i присвоена цена $p^{(j)}$, и нулю в противном случае. В предположении, что значение Q задано, задачу оптимизации можно определить как следующую целочисленную программу:

$$\begin{aligned} & \max_z \quad \sum_{i \in N} \sum_{j \in P} p^{(j)} \cdot q_i(p^{(j)}, Q) \cdot z_{ij} \\ & \text{при условии} \quad \sum_{j \in P} z_{ij} = 1 \\ & \quad \sum_{i \in N} \sum_{j \in P} p^{(j)} \cdot z_{ij} = Q \\ & \quad z_{ij} \in \{0, 1\}. \end{aligned} \quad (6.95)$$

Первое ограничение гарантирует, что каждый продукт имеет в точности одну цену, а второе — что сумма всех цен на продукты равна Q . Функция спроса q_i про-

гнозирует спрос на продукт i как функция цены, присвоенной продукту, и суммы цен Q , то есть сумма цен используется в обучении модели как один из признаков. Функция спроса, как уже говорилось, может учитывать доступный запас, поэтому прогнозируемый уровень спроса можно ограничить уровнем запаса продукта. Это особенно важно, если модель используется для планирования распродаж и оптимизации цен, когда целью является продажа запасов.

Задача целочисленного программирования 6.95 имеет существенно меньшую вычислительную сложность, чем упрощенный подход полного перебора вариантов, но она все еще может быть слишком сложной при большом количестве продуктов и цен. В таких случаях можно использовать ослабленную задачу линейного программирования (где z_{ij} — не бинарные, а непрерывные переменные, такие, что $0 \leq z_{ij} \leq 1$), чтобы примерно оценить все возможные значения Q и затем найти точное оптимальное решение путем решения задачи целочисленного программирования 6.95 для подмножества из Q значений [Ferreira et al., 2016]. Такой подход позволяет еще больше снизить вычислительную сложность и сделать задачу пригодной для практического применения даже для большого количества продуктов и ценовых уровней.

6.7.3. Персонализированные скидки

Методы оптимизации, которые мы рассмотрели выше, используют разные степени готовности платить, устанавливая различные скидки для различных клиентских сегментов или временных интервалов. В конечном итоге хотелось бы объединить эти два подхода для управления денежными и временными свойствами скидок на уровне сегмента. Кроме того, можно попытаться повысить эффективность дифференциации цен, заменив сегментацию с персонализированными скидками. На этом этапе методы ценообразования сходятся с методами продвижения, описанными в главе 3: службы ценообразования могут использовать преимущества методов таргетирования для принятия решений о ценообразовании на основе индивидуальных профилей клиентов, а службы продвижения — оптимизировать денежные аспекты продвижения, такие как величина скидки, с помощью методов оптимизации цен. Далее в этом разделе мы разработаем метод, оптимизирующий величину скидки и пытающийся найти оптимальное время и продолжительность для предложения скидки данному пользователю [Johnson et al., 2013]. Идея оптимизации временных свойств обусловлена предположением, что вероятность покупки клиентом неоднородна и меняется с течением времени, поэтому для каждого пользователя существует оптимальное временное окно скидки.

Чтобы смоделировать временные свойства скидки, разложим вероятность покупки бренда k покупателем u в момент времени t со значением скидки d на два множи-

теля: вероятность покупки бренда и вероятность совершить покупку в момент времени t :

$$p_{k|t|d} = p(\text{brand} = k | u; d) \cdot p(\text{time} = t | u; d). \quad (6.96)$$

Теперь смоделируем функции плотности вероятности $p(\text{brand} = k | u; d)$ и $p(\text{time} = t | u; d)$ по отдельности, но для обеих используем общий подход. Сначала определим форму распределения вероятностей и опишем ее с помощью функции полезности в качестве параметра. Затем построим регрессионную модель для оценки функции полезности на основе данных.

Функция плотности вероятности покупки данного бренда $p(\text{brand} = k | u; d)$ является типичным случаем модели множественного выбора, когда потребитель выбирает бренд из нескольких взаимозаменяемых альтернатив (обозначим общее количество конкурирующих брендов как K). Следовательно, для определения распределения можно использовать полиномиальную логит-модель (Multinomial Logit, MNL):

$$p(\text{brand} = k | u; d) = \frac{\exp(x_{kut})}{\sum_{i=1}^K \exp(x_{iut})}. \quad (6.97)$$

Функцию полезности x_{kut} можно определить по данным, построив, например, такую регрессионную модель:

$$x_{kut} = \sum_{w=1}^W \beta_{uw} F_{kutw}, \quad (6.98)$$

где F_{kutw} — W объясняющих переменных, в число которых входит скидка d и другие признаки, такие как лояльность и цена, а β_{uw} — коэффициенты регрессии W .

Функция плотности вероятности покупки в момент времени t моделируется в [Johnson et al., 2013] в виде распределения Эрланга:

$$p(\text{time} = t | u; d) = y_u^2 \cdot t \cdot \exp(-y_u t), \quad (6.99)$$

где переменный параметр y_u можно определить с помощью регрессионной модели, которая, подобно модели для переменного параметра x в уравнении 6.97, включает в число объясняющих переменных значение скидки и в дальнейшем может быть предметом оптимизации.

Определенная выше вероятность покупки позволяет смоделировать объем продаж для данного клиента Q_u как функцию значения скидки в долларах d , времени начала t и продолжительности действия скидки T :

$$Q_u(d, t, T) = \int_t^{t+T} p_{klud} dt. \quad (6.100)$$

В результате мы приходим к следующей задаче оптимизации валовой прибыли:

$$\max_{d, t, T} \sum_u m \cdot (Q_u(0, 0, t) + Q_u(d, t, t+T) + Q_u(0, t+T, \infty)) - d \cdot Q_u(d, t, T), \quad (6.101)$$

где m — маржа для обычной цены. Первый член в уравнении выше соответствует доходу, который, в свою очередь, состоит из трех компонентов — доход, полученный до начала, во время и после акции — и второй член соответствует затратам на рекламу. Это деление иллюстрирует диаграмма на рис. 6.19.

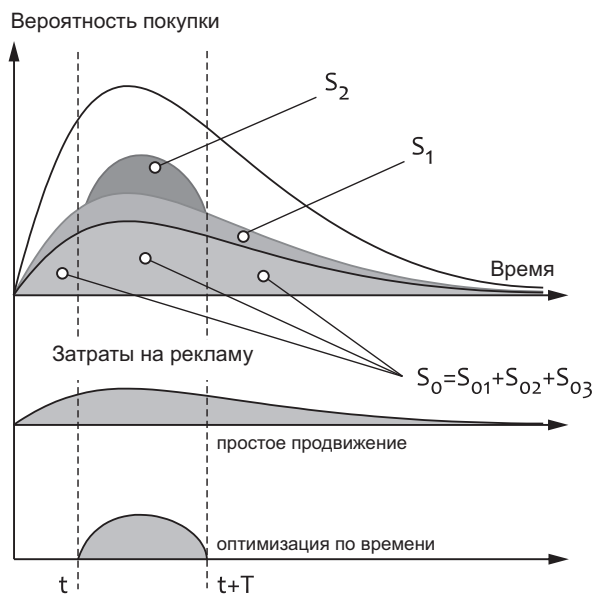


Рис. 6.19. Оптимизация временного окна для продвижения

График в верхней части показывает плотность вероятности покупки клиентом u , когда ожидаемый объем продаж для данного продукта по обычной цене соответствует области S_0 . Плоская постоянная скидка увеличивает этот объем, добавляя область S_1 , поэтому общий доход и рекламные расходы (показанные на среднем графике) будут пропорциональны $S_0 + S_1$. Продвижение, оптимизированное по времени, делает доход пропорциональным сумме $S_0 + S_2$, а затраты на него будут пропорциональны $S_{02} + S_2$ (график внизу). Эта разница между плоским и опти-

мизированным продвижением показывает потенциал использования временной оптимизации в случае определенных количественных свойств функций плотности вероятности.

6.8. Распределение ресурсов

Динамическое ценообразование предлагает способ сегментирования клиентов по их готовности платить и оптимизировать цены для каждого сегмента с учетом ограниченной общей емкости. Как уже говорилось выше, одно из самых больших преимуществ такого подхода заключается в том, что он не требует определять ценовые сегменты заранее и может создавать и настраивать их динамически. С другой стороны, он требует, чтобы бизнес-модель и среда были достаточно гибкими в установлении и обновлении цен. Эта гибкость во многом зависит от отрасли. Например, в розничной торговле и электронной коммерции обычно имеются хорошие возможности для динамического ценообразования, тогда как отрасли услуг, такие как авиакомпании, гостиницы и грузовые перевозки, могут быть менее гибкими в этом отношении. Это различие отчасти объясняется историческими причинами, а именно практикой установления фиксированных тарифов для различных классов обслуживания. Это привело к разработке большой группы методов, основанных на альтернативной интерпретации задачи. Эти методы впервые появились в отрасли авиаперевозок и исторически предшествовали появлению методов динамического ценообразования, а также большинства программных методов в целом.

Предположив, что продавец имеет операционные, юридические или бизнес-ограничения, влияющие на возможность произвольно менять цены, задачу динамической оптимизации цен можно перевернуть с ног на голову и рассмотреть альтернативный подход. Идея состоит в том, чтобы определить набор фиксированных ценовых сегментов, обычно называемых *классами тарифов*, и распределить часть общей емкости между классами так, чтобы максимизировать прибыль. Соответственно, предметом оптимизации становятся лимиты емкости, выделенные для каждого класса. Классическим примером этой задачи может служить авиакомпания, предлагающая три класса тарифов (например, эконом, бизнес и первый класс) и решающая, сколько мест зарезервировать для каждого класса с учетом фиксированной общей вместимости самолета.

6.8.1. Среда

Выше мы установили, что динамическое ценообразование возможно для сред с определенными свойствами, такими как изменчивость спроса и фиксированная емкость ресурсов. Эти фундаментальные соображения в целом применимы

и к распределению ресурсов, поскольку, по сути, это всего лишь другое решение той же задачи. Однако многие методы распределения ресурсов были разработаны главным образом для сферы услуг и устраняют ряд характерных для нее ограничений. Далее мы ограничимся изучением базовой среды и лишь кратко рассмотрим дополнительные сложности, существующие в теории и практике распределения ресурсов.

- Продавец предлагает товар или услугу нескольким сегментам рынка по разным тарифам. Сегменты могут определяться уровнями обслуживания, такими как экономический или бизнес-класс в авиационных перевозках, или основываться на более сложных бизнес-правилах, направленных на более совершенное разграничение классов. Например, сеть отелей или авиакомпания может быть готова продавать свои услуги бизнес-клиентам по более высокой цене, чем рядовым отдыхающим. Чтобы бизнес-клиенты не могли приобрести услугу по более низкой цене, поставщик услуг может установить условие, согласно которому предложения по низкой цене должны бронироваться за несколько недель или исключать возможность ночевки в субботу.
- Все тарифные классы обслуживаются из одной и той же фиксированной емкости ресурса, но лимиты бронирования для классов могут изменяться динамически. Например, авиакомпания может выделить на разные рейсы разный процент посадочных мест стандартного экономического и льготного классов, но оставить общее количество мест экономического класса неизменным.
- Система оптимизации анализирует исторические и текущие данные о спросе и по результатам анализа определяет или обновляет лимиты бронирования для каждого класса. Лимиты загружаются в систему бронирования, которая принимает запросы на бронирование. Запрос на бронирование — это запрос на резервирование единицы емкости для указанного класса тарифа, например запрос на бронирование одного места эконом-класса со скидкой. Система бронирования либо принимает запрос, если соответствующий лимит больше нуля, с последующим уменьшением счетчика лимита, либо отклоняет запрос, если емкость исчерпана.
- Базовые модели оптимизации, которые мы рассмотрим ниже, делают несколько важных предположений о спросе. Во-первых, предполагается, что спрос на каждый класс является случайной величиной с известным распределением. Во-вторых, предполагается, что все переменные спроса независимы. В частности, спрос на данный класс не зависит от наличия или отсутствия других классов. Это очень грубое приближение, потому что клиент с запросом на определенный тарифный класс может рассмотреть другие классы в случае отклонения заявки и увеличить соответствующую группу спроса, в точности

как в других случаях с несовершенной сегментацией. Наконец, предполагается, что запросы поступают последовательно от самого низкого класса (самого дешевого) к самому высокому (самому дорогому). Это предположение тоже является относительно грубым приближением, но оно широко используется в практике и часто соответствует реальным моделям спроса — например, правила ограждения бюджетных туристов часто включают условие, требующее произвести бронирование заранее.

Следует отметить, что многие приложения распределения ресурсов требуют выполнения двух основных условий, не включенных в нашу базовую модель среды. Во-первых, ресурсы часто выделяются не как отдельные единицы, а как продукты, включающие несколько единиц. Например, бронирование отеля — это продукт, включающий одну или несколько ночей проживания, а маршрут перелета может состоять из цепочки рейсов. Это требует единого управления и оптимизации сети ресурсов. Во-вторых, во многих отраслях, включая авиакomпаний и отели, запросы на бронирование могут отменяться, и доля отмененных запросов может быть значительной. Например, авиакомпания American Airlines сообщает, что около половины бронирований отменяются или пассажиры не являются на рейс [Smith et al., 1992]. Это приводит к практике *избыточного бронирования*, когда поставщик услуг позволяет резервировать ресурсы за пределами емкости, ожидая, что некоторые из этих резервирований будут отменены в будущем. Избыточное бронирование также требует разработки специализированных методов, корректирующих уровни бронирования в соответствии с ожидаемым количеством отмен.

Как уже упоминалось, задача распределения ресурсов рассматривает лимиты бронирования как предмет оптимизации. Самый простой способ определить лимиты бронирования — выделить определенную емкость отдельно для каждого тарифного класса. Основная проблема этого подхода заключается в том, что емкость для более высокого класса может оказаться исчерпана, тогда как более низкие классы останутся доступными. В результате система бронирования будет отклонять запросы на более высокий класс, чтобы сохранить емкость, которую она сможет использовать для будущих запросов на низкий тариф. Такое поведение, очевидно, наносит ущерб с точки зрения прибыльности. Лучшее решение, принятое в качестве стандарта в большинстве теоретических моделей и практических приложений, предлагает подход на основе *вложенных пределов*. Идея состоит в том, чтобы установить ограничения не для конкретного тарифного класса, а для всех классов выше или равных данному, как показано на рис. 6.20. Эти пределы, их часто называют *уровнями защиты*, устанавливаются так, что первый предел равен емкости, зарезервированной для первого класса, а последний предел — общей емкости.

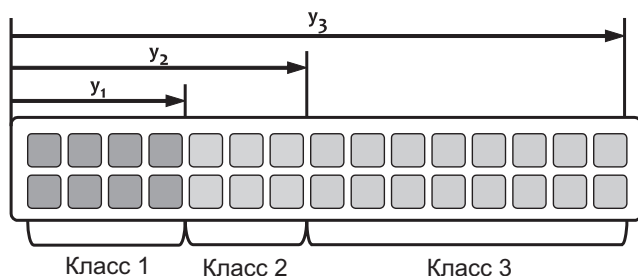


Рис. 6.20. Три тарифных класса и уровни защиты y_1 , y_2 и y_3

Система бронирования принимает или отклоняет запрос в соответствии со следующей логикой:

- Запрос на класс i принимается, только если $y_i > y_{i-1}$. Это очевидно, потому что предел резервирования для класса i — это разность между y_i и y_{i-1} .
- Если запрос принят, y_n уменьшается, и все уровни защиты, превышающие новое значение y_n , устанавливаются равными y_n . Это можно представить на рис. 6.20 как правую границу y_n , перемещаемую влево по мере резервирования все большей и большей доли емкости; правые границы других уровней фиксируются.

Такой процесс бронирования показан в табл. 6.8. Изначально у нас есть 2 единицы первого класса, 2 единицы второго класса и 4 единицы третьего класса, что в сумме дает общую емкость, равную 8. До 5-го раунда принимается любой запрос, когда общая емкость уменьшится до 4 единиц. Разность между y_3 и y_2 равна нулю, поэтому класс 3 закрывается. Обратите внимание, что этот процесс необратим: как только разность между y_i и y_{i-1} становится равной нулю, класс i остается постоянно закрытым. Заявки на второй класс принимаются до 7 раунда, после чего принимаются только заявки на первый класс.

Таблица 6.8. Пример процесса бронирования с вложенными пределами

	Уровни защиты			Продано единиц			Запрос на бронирование	Действие
	y_1	y_2	y_3	C_1	C_2	C_3		
1	2	4	8	0	0	0	1 единица в классе C_2	Принято
2	2	4	7	0	1	0	1 единица в классе C_3	Принято
3	2	4	6	0	1	1	1 единица в классе C_3	Принято

	Уровни защиты			Продано единиц			Запрос на бронирование	Действие
	y_1	y_2	y_3	C_1	C_2	C_3		
4	2	4	5	0	1	2	1 единица в классе C_1	Принято
5	2	4	4	1	1	2	1 единица в классе C_3	Отклонено
6	2	4	4	1	1	2	1 единица в классе C_1	Принято
7	2	3	3	2	1	2	1 единица в классе C_2	Принято
8	2	2	2	2	2	2	1 единица в классе C_2	Отклонено
9	2	2	2	2	2	2	1 единица в классе C_3	Отклонено
10	2	2	2	2	2	2	1 единица в классе C_1	Принято
11	1	1	1	3	2	2	1 единица в классе C_1	Принято
12	0	0	0	4	2	2	–	Отклонено

6.8.2. Распределение с двумя классами

Распределение — сложная задача, поэтому начнем с самого простого сценария — с двумя классами тарифов. Предположим, что имеется емкость C единиц, и обозначим цены для первого и второго классов как p_1 и p_2 соответственно, где $p_1 > p_2$.

Согласно нашим предположениям о среде, спрос на каждый класс является случайной величиной Q_i и известна ее кумулятивная функция распределения F_i . Запросы на выделение поступают последовательно, и, согласно другому предположению, сначала поступают запросы на второй, менее дорогой класс. Соответственно, наша цель — определить такой оптимальный уровень защиты y , чтобы принять не больше $C - y$ запросов на второй класс и оставшиеся y единиц зарезервировать для удовлетворения запросов на первый класс.

Каждый раз, когда поступает запрос на второй класс, мы можем принять его или отклонить и переключиться на пространство первого класса. Это решение легко проанализировать с точки зрения ожидаемых результатов, как показано на рис. 6.21. Принимая запрос, мы получаем доход p_2 . Но отклоняя его, мы закрываем второй класс и переключаемся на первый и оказываемся перед двумя вариантами развития событий. В первом случае, если спрос на первый класс превысит оставшуюся емкость y , мы продадим ее по цене p_1 . Во втором случае, если спрос на первый класс окажется ниже оставшейся емкости, часть ресурсов не будет забронирована

вообще, и мы получим нулевой доход. Следовательно, условие принятия запросов на второй класс можно записать следующим образом:

$$p_2 \geq p_1 \cdot \Pr(Q_1 \geq y), \quad (6.102)$$

что то же самое, что и

$$p_2 \geq p_1 \cdot (1 - F_1(y)). \quad (6.103)$$

Инвертируя кумулятивную функцию распределения, найдем оптимальный уровень защиты для первого класса:

$$y_{opt} = F_1^{-1} \left(1 - \frac{p_2}{p_1} \right). \quad (6.104)$$

Уравнение 6.104 известно как правило Литтлвуда [Littlewood, 1972]. Оптимальный уровень защиты не зависит от распределения спроса на второй класс, потому что мы пытаемся определить наибольший объем, который можно зарезервировать для первого класса, и предполагаем, что все остальное будет забронировано во втором классе, независимо от распределения.

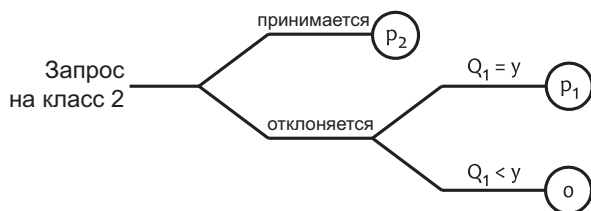


Рис. 6.21. Дерево решений для задачи распределения с двумя классами

ПРИМЕР 6.7

Проиллюстрируем задачу оптимизации для тарифов с двумя классами на следующем примере. Рассмотрим поставщика услуг, который имеет 20 единиц ресурса и продает их по полной цене 300 и сниженной цене 200 долларов. Спрос на услуги по полной цене, согласно оценкам, имеет нормальное распределение со средним значением 8 и стандартным отклонением 2. Следова-

тельно, вероятность, что спрос на услугу по полной цене превысит y единиц, можно выразить с помощью функции кумулятивного распределения Φ для стандартного нормального распределения:

$$\Pr(Q_1 \geq y) = 1 - \Phi\left(\frac{y - 8 - 0,5}{2}\right). \quad (6.105)$$

Обратите внимание, что из-за дискретной природы резервирования мы добавили сдвиг 0,5 — в этом случае вероятность, что спрос составит ровно y единиц, можно аппроксимировать интегрированием кумулятивной функции распределения на интервале от $y - 0,5$ до $y + 0,5$. Учитывая, что мы резервируем y единиц для продажи по полной цене, предельный доход для единицы по полной цене можно определить как

$$r_1(y) = \$300 \times \Pr(Q_1 \geq y). \quad (6.106)$$

Иначе говоря, предельный доход — это разность между ожидаемым доходом в сегменте с полной ценой, для которого выделено y единиц, и соответствующим доходом, когда выделено только $y - 1$ единиц. Общий ожидаемый доход в сегменте с полной ценой равен сумме этих предельных доходов:

$$R_1(y) = \sum_{i=1}^y r_1(i). \quad (6.107)$$

Доход в сегменте с льготной ценой — это просто количество оставшихся единиц, умноженное на льготную цену:

$$R_2(y) = \$200 \times (C - y). \quad (6.108)$$

Общий доход, полученный поставщиком, определяется как сумма доходов в сегментах с полной и льготной ценой. Теперь рассчитаем эти метрики для всех возможных значений уровня защиты y и поместим их в табл. 6.9.

Как видите, наибольший доход достигается при уровне защиты 7, то есть когда 7 единиц ресурса выделяются для сегмента с полной ценой и 13 единиц — для сегмента с льготной ценой. Правило Литтлвуда дает тот же результат — предельный доход $r_1(y)$, соответствующий правой части уравнения 6.102, опускается ниже льготной цены 200, начиная с уровня защиты 8.

Таблица 6.9. Пример оптимизации уровня защиты для тарифа с двумя классами

y	$r_1(y)$	$R_1(y)$	$R_2(y)$	$R_1(y) + R_2(y)$
1	299	299	3800	4099
2	299	599	3600	4199
3	299	898	3400	4298
4	296	1195	3200	4395
5	287	1483	3000	4483
6	268	1751	2800	4551
7	232	1983	2600	4583
8	179	2163	2400	4563
9	120	2283	2200	4483
10	67	2351	2000	4351
11	31	2383	1800	4183
12	12	2395	1600	3995
13	3	2398	1400	3798
14	0	2399	1200	3599
15	0	2399	1000	3399
16	0	2400	800	3200
17	0	2400	600	3000
18	0	2400	400	2800
19	0	2400	200	2600
20	0	2400	0	2400

6.8.3. Распределение с несколькими классами

Правило Литтлвуда предлагает компактное выражение для задачи распределения с двумя классами. На практике, однако, чаще возникает задача распределения с большим количеством классов. Найти оптимальное решение в этом случае намного сложнее, но есть несколько способов, способных помочь справиться с задачей. Один из них — использовать рекурсивный подход, исходя из предположения

о последовательном исчерпании классов спроса, когда выбор уровня защиты для одного класса уменьшает размерность задачи. Этот способ позволяет выразить и решить задачу распределения в терминах динамического программирования. Другой подход заключается в расширении вероятностного анализа, который мы применили для правила Литтлвуда, и использовании моделирования для поиска оптимальных уровней защиты [Brumelle and McGill, 1993; Talluri and Van Ryzin, 2004]. Далее мы рассмотрим этот последний подход.

Как мы уже видели, оптимальный уровень защиты для первого класса в задаче с двумя классами определяется по формуле

$$p_2 = p_1 \cdot \Pr(Q_1 \geq y_1^{opt}). \quad (6.109)$$

Давайте сделаем еще шаг и рассмотрим дерево решений для третьего класса, как показано на рис. 6.22.

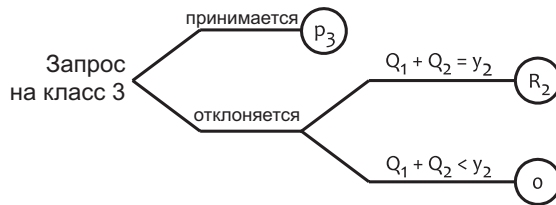


Рис. 6.22. Дерево решений для задачи распределения с тремя классами

Аналогично задаче с двумя классами запрос на третий класс может быть принят, что принесет нам доход p_3 , или отклонен. Последнее означает, что третий класс будет закрыт и все последующие запросы будут обрабатываться в режиме с двумя классами. Это приводит нас к двум возможным результатам:

- Если суммарный спрос на первый и второй классы окажется ниже уровня защиты y_2 , часть ресурсов не будет распределена — мы закрыли третий класс слишком рано.
- Иначе остальные запросы будут обрабатываться, как в стандартной задаче с двумя классами. Как было показано в предыдущем разделе, средний доход на единицу R_2 в этом случае равен p_2 , если выбран оптимальный уровень защиты y_1 , то есть в соответствии с правилом Литтлвуда.

Следовательно, оптимальное значение для y_2 можно выразить так:

$$p_3 = p_2 \cdot \Pr(Q_1 + Q_2 \geq y_2^{opt} | Q_1 \geq y_1^{opt}). \quad (6.110)$$

Мы можем сравнить уравнения 6.109 и 6.110 и рекурсивно применить подход на основе дерева решений, чтобы найти следующий оптимальный уровень защиты:

$$\frac{p_{j+1}}{p_j} = \Pr(Q_1 + \dots + Q_j \geq y_j^{opt} | Q_1 \geq y_1^{opt} \text{ AND } \dots \text{ AND } Q_1 + \dots + Q_{j-1} \geq y_{j-1}^{opt}). \quad (6.111)$$

Несмотря на то что уравнение 6.111 выглядит сложнее, чем правило Литтлвуда, оно дает относительно простой способ оценки уровня защиты путем моделирования. Для удобства рассмотрим пример с тремя классами, хотя этот метод легко можно распространить на случай с любым числом классов. Предположим, что распределения спроса Q_1 и Q_2 известны, и мы можем сгенерировать относительно большое количество двумерных точек с координатами, которые определяются как Q_1 и $Q_1 + Q_2$.

Оптимальное значение уровня защиты y_1 можно оценить с помощью уравнения 6.109 — мы должны найти линию, которая разбивает точки по первой координате так, чтобы количество точек слева и справа были разделены в той же пропорции, что и p_1 и p_2 , как показано на рис. 6.23. Точки справа удовлетворяют условию $Q_1 \geq y_1^{opt}$ из уравнения 6.110, поэтому, чтобы оценить уровень защиты y_2 , разделим их по второй координате так, чтобы количество точек в нижней и верхней части были разделены в той же пропорции, что и p_3 и p_2 .

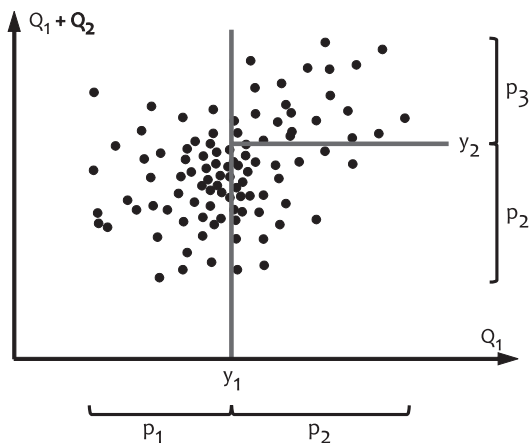


Рис. 6.23. Оптимизация уровней защиты для трех классов методом моделирования

6.8.4. Эвристические решения для нескольких классов

Несмотря на то что оптимальные уровни защиты можно рассчитать с использованием метода моделирования, рассмотренного в предыдущем разделе, а также альтернативных алгоритмов, многие реализации используют более простые эв-

ристические методы, дающие неоптимальные, но часто близкие к оптимальным решения. Наиболее важным методом такого рода является алгоритм расчета ожидаемого дохода на место (Expected Marginal Seat Revenue, EMSR), который имеет две версии: EMSRa и EMSRb [Belobaba, 1987, 1989]. Обе версии пытаются адаптировать правило Литтлвуда к случаю с несколькими классами.

6.8.4.1. EMSRa

Напомним, что уровень защиты для класса j — это общая емкость, зарезервированная для всех классов от j до 1. После расчета уровня защиты для более дешевых классов от n до $j + 1$ уровень защиты для класса j определяет, как разделить емкость между классом $j + 1$ и более дорогими классами. Идея алгоритма EMSRa состоит в том, чтобы аппроксимировать этот уровень защиты суммой полученных уровней, путем отдельного применения правила Литтлвуда к классу $j + 1$ и каждому из классов от j до 1. То есть сначала вычисляется j попарных уровней защиты с использованием следующих уравнений:

$$\begin{aligned} p_{j+1} &= p_j \Pr(Q_j \geq y_{j+1}^{(j)}) \\ p_{j+1} &= p_{j-1} \Pr(Q_{j-1} \geq y_{j+1}^{(j-1)}) \\ &\vdots \\ p_{j+1} &= p_1 \Pr(Q_1 \geq y_{j+1}^{(1)}). \end{aligned} \quad (6.112)$$

Окончательный уровень защиты вычисляется суммированием попарных уровней, как показано на рис. 6.24.

$$y_j = \sum_{k=1}^j y_{j+1}^{(k)}. \quad (6.113)$$

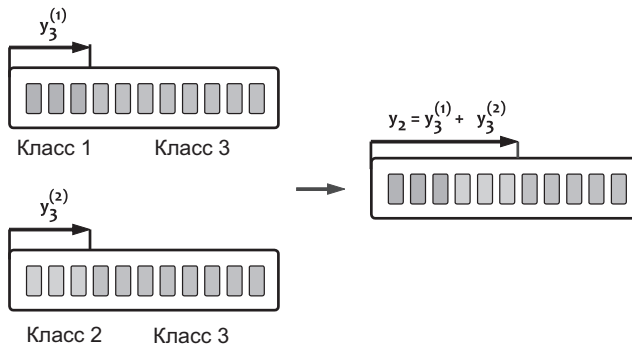


Рис. 6.24. Пример алгоритма EMSRa для трех классов

Сравнив уравнения 6.112 и 6.113 с оптимальным решением 6.111, можно заметить, что для аппроксимации вероятности превышения суммарным спросом определенного уровня алгоритм EMSRa использует вероятности превышения отдельными запросами уровней, определяемых соответствующими ценовыми отношениями. То есть алгоритм EMSRa дает чрезмерно консервативную оценку, в том смысле что резервирует слишком много единиц для более высоких классов и отклоняет слишком много запросов по низким тарифам.

6.8.4.2. EMSRb

Альтернативный подход к решению задачи заключается в слиянии классов от j до 1 в один виртуальный агрегатный класс, имеющий свой спрос и цену, с последующим применением правила Литтлвуда. Спрос на агрегированный класс можно оценить как сумму запросов для включенных в него классов:

$$\bar{Q}_j = \sum_{k=1}^j Q_k. \quad (6.114)$$

«Цену» агрегированного класса можно определить как средневзвешенную цену включенных классов:

$$\bar{p}_j = \frac{\sum_{k=1}^j p_k \mathbb{E}[Q_k]}{\sum_{k=1}^j \mathbb{E}[Q_k]}. \quad (6.115)$$

Уровень защиты для класса j можно оценить, применив правило Литтлвуда к агрегированному классу и предыдущему классу $j + 1$, то есть

$$p_{j+1} = \hat{p}_j \cdot \Pr(\hat{Q}_j \geq y_j). \quad (6.116)$$

По общему признанию, оба алгоритма, EMSRa и EMSRb, дают аппроксимации, очень близкие к оптимальным решениям. Алгоритм EMSRb разрабатывался как улучшенная версия EMSRa и по некоторым отзывам иногда дает лучшие результаты, чем EMSRa, но по итогам экспериментов с реальными данными ни один из методов не дает неизменно лучших результатов [Talluri and Van Ryzin, 2004].

6.9. Оптимизация ассортимента

На готовность клиента платить за данный продукт или услугу почти всегда влияют альтернативные варианты, то есть возможность выбора конкурирующего или заменяющего предложения. Влияние этих сил на функции спроса и, в конечном

итоге, на прибыль может быть более или менее значительным, в зависимости от конкретной отрасли и ситуации. Мы уже знаем, что в случае смещения или зависимо-го спроса в рамках категории ценовые решения для нескольких продуктов или клиентских сегментов должны оптимизироваться совместно. Эта задача становится еще более сложной в розничной торговле, где приходится иметь дело с десятками и сотнями тысяч продуктов, поэтому не только цены, но и другие маркетинговые и корпоративные ресурсы должны оптимизироваться с учетом зависимостей спроса и эффектов замещения. В этом разделе мы подробно рассмотрим данный тип задач оптимизации.

6.9.1. Оптимизация планировки магазина

Зависимости спроса между продуктами и категориями и их сходство предполагают возможность перекрестной продажи связанных или дополнительных предложений при приобретении покупателями определенных продуктов. Примерами таких дополнительных продуктов могут служить кофе и сахар, косметика и сумки, и даже пиво и подгузники. Один из способов использования этих возможностей — определение продуктов, которые клиенты склонны покупать вместе, и оптимизация планировки магазина или размещения контента на сайте, чтобы упростить навигацию и подтолкнуть клиентов к покупке наборов вместо отдельных продуктов.

Первую задачу этой стратегии — определение продуктов, которые часто покупают-ся вместе, — можно решить с помощью простого анализа покупательской корзины. Предположим, что у нас есть история сделок, где каждая сделка t представлена набором элементов r , купленных в рамках этой сделки:

$$t_n = \{r_{n1}, r_{n2}, \dots, r_{nk}\}. \quad (6.117)$$

Поддержка (support) элемента или набора элементов — это доля сделок в истории, содержащих элемент или набор элементов. Иначе говоря, это эмпирическая веро-ятность, что случайно выбранная сделка содержит данный элемент или набор эле-ментов. *Подъем* (lift) определяется для пары элементов как отношение поддержки для пары к произведению поддержек элементов по отдельности:

$$\lambda(r_a, r_b) = \frac{\text{support}(r_a \text{ AND } r_b)}{\text{support}(r_a) \times \text{support}(r_b)}. \quad (6.118)$$

Другими словами, подъем — это отношение наблюдаемой вероятности со-вместного появления двух элементов к совместной вероятности появления, рассчитанной в предположении их независимости. То есть подъем выше еди-

ницы указывает на сходство элементов (разумеется, с учетом статистической значимости результатов). Подъем можно измерить не только для пар продуктов, но и для пар категорий, сопоставив каждый элемент в истории сделок со своей категорией и оценив выражение 6.118, исходя из предположения, что r_a и r_b являются категориями.

Теперь вернемся к задаче оптимизации планировки магазина. У нас есть n категорий продуктов и n местоположений, таких как проходы или полки. Задачу оптимизации в таком случае можно сформулировать как размещение всех категорий в разных местоположениях такое, что категории с высоким попарным сходством окажутся близко друг к другу. Для начала рассчитаем матрицу попарных подъемов для категорий:

$$L = \{\lambda_{ij}\}, \quad i, j = 1, \dots, n. \quad (6.119)$$

Далее нам потребуется матрица расстояний между местоположениями:

$$D = \{d_{ij}\}, \quad i, j = 1, \dots, n, \quad (6.120)$$

где d_{ij} — расстояние между местоположениями i и j . Расстояние может определяться разными способами, например, это может быть бинарная переменная, равная единице, если местоположения находятся рядом, и нулю в противном случае. Теперь задачу оптимизации можно определить как

$$\max_{\pi} \sum_{i=1}^n \sum_{j=1}^n \lambda_{ij} d_{\pi(i), \pi(j)}, \quad (6.121)$$

где $\pi(i)$ — функция перестановки, отображающая категории в местоположении, то есть $\pi(x)$ равна y , когда категория с номером x присваивается местоположению y . Задача 6.121 является примером квадратичной задачи о назначениях¹ (Quadratic Assignment Problem, QAP). Это хорошо изученная комбинаторная задача оптимизации, которая, однако, имеет высокую вычислительную сложность. Тем не менее есть сведения, что этот метод использовался на практике для оптимизации планировки в небольших магазинах [Winston, 2014].

¹ Впервые задача QAP была представлена в контексте исследования операций для моделирования следующей задачи. Есть множество предприятий и местоположений. Цель состоит в том, чтобы разместить предприятия в этих местоположениях так, чтобы минимизировать общую стоимость, где стоимость назначения определяется как произведение суммы расстояний между местоположениями на соответствующие потоки между предприятиями.

ПРИМЕР 6.8

Следующий пример демонстрирует решение задачи оптимизации планировки магазина. Рассмотрим продуктовый магазин, предлагающий шесть категорий продуктов: молочные продукты, деликатесы, хлебобулочные изделия, напитки, овощи и замороженные продукты. Матрица попарных подъемов для этих категорий оценивается по историческим данным:

$$L = \begin{array}{l} \text{Молоко} \\ \text{Деликатесы} \\ \text{Хлеб} \\ \text{Напитки} \\ \text{Овощи} \\ \text{Заморозка} \end{array} \begin{array}{c} \begin{array}{cccccc} \text{Молоко} & \text{Деликатесы} & \text{Хлеб} & \text{Напитки} & \text{Овощи} & \text{Заморозка} \end{array} \\ \left[\begin{array}{cccccc} 1,00 & 0,80 & 1,30 & 0,90 & 1,00 & 0,90 \\ 0,80 & 1,00 & 1,20 & 1,10 & 1,30 & 0,80 \\ 1,30 & 1,20 & 1,00 & 1,30 & 1,20 & 0,90 \\ 0,90 & 1,10 & 1,30 & 1,00 & 1,20 & 1,50 \\ 1,00 & 1,30 & 1,20 & 1,20 & 1,00 & 0,80 \\ 0,90 & 0,80 & 0,90 & 1,50 & 0,80 & 1,00 \end{array} \right] \end{array} \quad (6.122)$$

На рис. 6.25 показан план магазина — сетка 2×3 , где каждое местоположение — это полка с товаром. В общем, у нас есть шесть свободных мест для шести категорий.

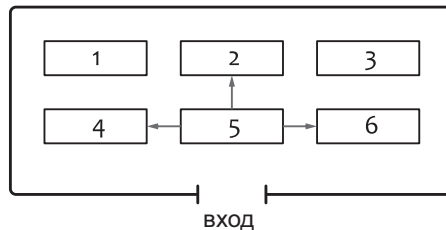


Рис. 6.25. План магазина

Матрица расстояний для этой планировки определяется следующим образом:

$$D = \begin{array}{l} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{array} \begin{array}{c} \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \end{array} \\ \left[\begin{array}{cccccc} 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{array} \right] \end{array} \quad (6.123)$$

Здесь предполагается, что расстояние равно единице только между соседними местами, иначе оно равно нулю. Например, место номер пять имеет расстояние, равное единице, до мест с номерами два, четыре и шесть. Решая задачу оптимизации 6.121 для указанных выше матриц, находим оптимальную планировку, представленную на рис. 6.26. Этот небольшой пример легко решить, оценив все $6! = 720$ возможных перестановок, но более крупные задачи требуют использования специализированного программного обеспечения, способного решить задачу QAP или одну из ее ослабленных версий.

6.9.2. Управление категориями

Задача управления категориями возникает, когда продавцу необходимо оптимизировать продажи всей товарной категории, а не отдельного товара. Эта задача очень типична для розничной торговли, потому что продавец может относительно легко изменить ассортимент в категории, и его основная задача — наиболее эффективно использовать имеющиеся ресурсы, такие как полки.

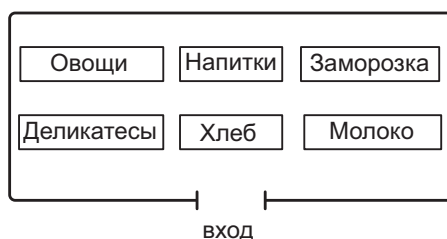


Рис. 6.26. Одна из оптимальных планировок магазина для случая с шестью категориями.

Конечно, альтернативные оптимальные планировки можно получить, отразив план по вертикали или по горизонтали

Категория — это набор относительно тесно связанных продуктов, имеющих много общего, таких как «молочные десерты» или «женские джинсы», поэтому, как правило, клиенты готовы заменить один продукт другим, если продукт по их выбору недоступен по какой-то причине. Недоступность может быть обусловлена преднамеренными изменениями ассортимента или исчерпанием складских запасов. Одна из главных целей управления категориями — найти подмножество продуктов, отвечающее физическим ограничениям, таким как доступное пространство на полке, и максимизировать прибыль за счет оптимального использования эффекта замещения. Также эту задачу можно сформулировать, как выявление наименее значимых продуктов, которые можно исключить из ассортимента и заместить

другими продуктами, не вызывая отрицательного влияния на прибыль. Результаты этого анализа можно применять при оптимизации различных аспектов:

- Уровни запасов продуктов можно оптимизировать, добавив учет эффектов замещения и потенциальных потерь, вызванных исчерпанием запасов.
- Планировку полок можно оптимизировать, чтобы скорректировать относительные доли продуктов на полке.
- Ассортимент можно оптимизировать путем добавления или удаления продуктов из ассортимента.

С точки зрения эконометрики задача управления категориями вытекает из закона уменьшения отдачи или, точнее говоря, из того факта, что доходы и затраты по-разному зависят от размера категории. Дело в том, что покупательная способность потребителей в какой-то момент достигает предела, тогда как затраты продолжают расти из-за увеличения торговых площадей и других операционных издержек, как показано на рис. 6.27. Эта тенденция приводит к задаче оптимизации категорий. Это очень сложная **задача**, потому что требует моделирования целой категории с учетом взаимозависимостей между продуктами в ней. Однако, несмотря на эти трудности, в сети супермаркетов Albert Heijn в Нидерландах [Kök and Fisher, 2007] была разработана и использована практически осуществимая модель оптимизации ассортимента. И оставшуюся часть раздела мы посвятим изучению этого решения.

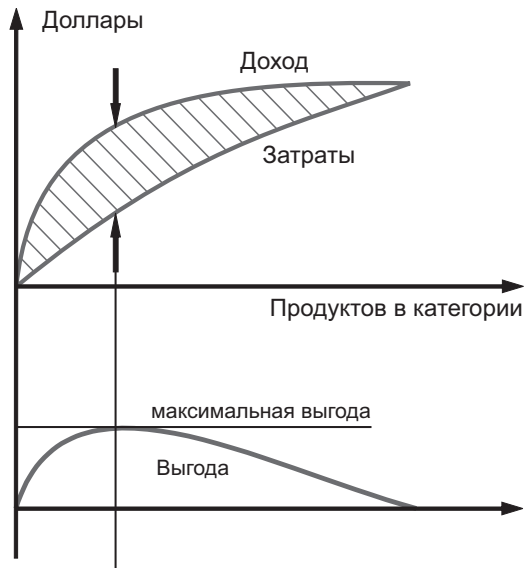


Рис. 6.27. Снижение отдачи в управлении категориями

Рассмотрим сеть супермаркетов, включающую несколько магазинов. Каждый магазин продает множество категорий товаров, но, как отмечалось выше, спрос считается зависимым только внутри категории, а сами категории считаются независимыми. Следовательно, ретейлер решает задачу оптимизации ассортимента для каждой категории в каждом магазине независимо. Для начала введем следующие обозначения, относящиеся к одной категории в одном магазине:

- $N = \{1, 2, \dots, J\}$: максимальное множество продуктов в категории, которую ретейлер предлагает своим клиентам, то есть полный ассортимент.
- $f_i \in \{0, 1, 2, \dots\}$: уровень запасов продукта j . Ретейлер оптимизирует свой ассортимент, выбирая f равным нулю (продукт отсутствует в ассортименте) или выше. Обозначим также вектор уровней запасов для всех J продуктов как $\mathbf{f} = \{f_1, \dots, f_J\}$.
- F_0 : общая емкость запасов, измеряется в тех же единицах, что и уровни запасов. Предполагается, что сумма уровней запасов для всех продуктов не может превышать F_0 . Общая емкость может быть ограничена складским пространством или доступным полочным пространством в магазине.
- $N_h \subset N$: ассортимент в магазине h , подмножество полного ассортимента.
- d_j : исходная норма спроса на товар j , то есть количество покупателей, которые выбрали бы товар j при наличии полного ассортимента N . Обозначим также вектор спроса на все товары как $\mathbf{d} = \{d_1, \dots, d_J\}$.
- D_j : наблюдаемый уровень спроса на продукты, то есть фактическое количество клиентов в день, выбравших продукт j согласно их первоначальному намерению или в результате эффекта замещения. Наблюдаемый спрос на продукт зависит от первичного спроса и наличия других продуктов, замещающих его, поэтому его можно рассматривать как функцию $D_j(\mathbf{f}, \mathbf{d})$.

Используя обозначения выше, задачу оптимизации ассортимента для данного магазина и данной категории можно выразить следующим образом:

$$\begin{aligned} & \max_{\mathbf{f}} \sum_{j \in N} G_j(f_j, D_j(\mathbf{f}, \mathbf{d})) \\ & \text{с условием } \sum_j f_j \leq F_0, \end{aligned} \tag{6.124}$$

где G_j — функция, описывающая прибыль для данного продукта и соответствующий наблюдаемый спрос. Эта функция сильно зависит от бизнес-модели ретейлера, и мы можем выделить несколько общих шаблонов, которые можно адаптировать для практического использования. Самый простой способ смоделировать прибыль — умножить наблюдаемый спрос на маржу продукта m :

$$G_j(f_j, D_j) = m_j \cdot D_j. \quad (6.125)$$

Уравнение 6.125 неявно предполагает идеальное пополнение запасов и отсутствие ситуаций их исчерпания. Это часто характерно для хорошо продаваемых потребительских товаров, таких как продукты питания, однако другие сферы розничной торговли, такие как торговля одеждой, тоже, вероятно, должны учитывать возможность исчерпания запасов, принимая во внимание минимальный спрос и фактический уровень запасов:

$$G_j(f_j, D_j) = m_j \cdot \min(D_j, f_j). \quad (6.126)$$

Ретейлеры, торгующие скоропортящимися товарами, должны также учитывать потери, связанные с ликвидацией запасов, которые можно смоделировать введением удельных потерь L , применяемых к непроданным запасам:

$$G_j(f_j, D_j) = m_j \cdot \min(D_j, f_j) - L_i \cdot (f_i - \min(D_j, f_j)). \quad (6.127)$$

В целях упрощения в дальнейшем будем считать, что все продукты своевременно пополняются, поэтому исчерпание складских запасов невозможно или незначительно. Это позволяет рассматривать $f_j \in \{0, 1\}$ как бинарную переменную, указывающую на наличие товара в ассортименте.

Для решения задачи оптимизации 6.124 необходимо определить функцию наблюдаемого спроса. В предположении неисчерпаемости запасов, которое мы сделали выше, функцию спроса можно смоделировать следующим образом:

$$D_j(\mathbf{f}, \mathbf{d}) = d_j + \sum_{k: f_k=0} \alpha_{k \rightarrow j} \cdot d_k, \quad (6.128)$$

где $\alpha_{k \rightarrow j}$ — вероятность замещения продукта k продуктом j . Формула 6.128 является относительно простой: первый член — это первичный спрос, а второй член соответствует кумулятивному эффекту замещения от всех продуктов, которые исключаются из ассортимента.

Уравнение 6.128 требует оценки вероятностей замещения $\alpha_{k \rightarrow j}$ и уровней первичного спроса d_j . Чтобы получить эти оценки, предположим, что известны следующие переменные (прогнозирование спроса мы уже рассматривали в разделе 6.6):

- Q_{jh} , $j \in N_h$: спрос на продукт j для одного клиента в магазине h . Если количество клиентов, посещающих магазин h в течение дня, обозначить как K_h , тогда $D_j = K_h \cdot Q_{jh}$.

- $Q_{jh}^0, j \in N$: спрос на продукт j для одного клиента в магазине h с полным ассортиментом (предположим, что существуют магазины с полным ассортиментом). Q_{jh}^0 соответствует первичному спросу, поскольку в магазинах с полным ассортиментом замещения не происходит.

Коэффициент замещения $\alpha_{k \rightarrow j}$ оценить труднее, потому что для ассортимента из J продуктов может существовать до J^2 разных вероятностей. Маловероятно, что у ретейлера имеется достаточно данных, чтобы надежно оценить это разнообразие. Однако есть эмпирические доказательства, что следующая упрощенная модель поведения клиента является достаточно точной на практике и требует оценки только одной переменной вместо J^2 : если продукт k недоступен, клиенты могут выбрать «запасной» продукт j с вероятностью δ , одинаковой для всех продуктов в категории, или отказываются от покупки с вероятностью $(1 - \delta)$. Эта модель приводит к следующей простой формуле определения коэффициента замещения:

$$\alpha_{k \rightarrow j} = \delta \frac{1}{|N|}. \quad (6.129)$$

Для оценки δ определим суммарный спрос в данном магазине как сумму значений Q_{jh} , которую можно определить по данным:

$$S_h \sum_{j \in N_h} Q_{jh}. \quad (6.130)$$

С другой стороны, это значение можно оценить по выражению 6.128:

$$\begin{aligned} \hat{S}_h(\delta) &= \sum_{j \in N_h} \left[Q_{jh}^0 + \sum_{k \in N \setminus N_h} \alpha_{k \rightarrow j} Q_{kh}^0 \right] = \\ &= \sum_{j \in N_h} Q_{jh}^0 + \sum_{j \in N_h} \sum_{k \in N \setminus N_h} \frac{\delta}{|N|} Q_{kh}^0. \end{aligned} \quad (6.131)$$

Теперь можно оценить δ , решив следующую задачу оптимизации, минимизирующую расхождение между наблюдаемыми и прогнозируемыми значениями суммарного спроса:

$$\delta_0 = \operatorname{argmax}_{0 \leq \delta \leq 1} \sum_h \left(\hat{S}_h(\delta) - S_h \right)^2. \quad (6.132)$$

Следующий шаг в решении задачи оптимизации 6.124 — вычисление норм первичного спроса, которые используются в уравнении 6.128. Сначала заметим, что общий спрос на все товары в N в магазине h можно вычислить следующим образом:

$$T_h = V_h \cdot \sum_{j \in N} Q_{jh}^0 \cdot \frac{S_h}{\hat{S}_h(\delta_0)}, \quad (6.133)$$

где V_h — общее количество клиентов, посещающих магазин h в течение дня. В уравнении 6.133 сумма всех Q_{jh}^0 , умноженная на V_h , представляет общий спрос при полном ассортименте. Однако значения Q_{jh}^0 оцениваются для магазинов с полным ассортиментом, поэтому специфика данного магазина h (например, местонахождение, размер магазина в квадратных метрах и т. д.) не моделируется. Это компенсируется масштабированием отношения оценочного спроса по категориям из уравнения 6.130 к прогнозируемому спросу из уравнения 6.131.

В магазине с ограниченным ассортиментом общий спрос T_h складывается из двух составляющих: спроса на продукты, входящие в ассортимент данного магазина, и спроса на другие продукты в N . Отношение между этими двумя компонентами можно выразить через Q_{jh}^0 :

$$r_h = \frac{\sum_{j \in N_h} Q_{jh}^0}{\sum_{j \in N} Q_{jh}^0}. \quad (6.134)$$

Соответственно, $T_h \cdot r_h$ представляет часть спроса, обусловленную продуктами в ассортименте, а $T_h \cdot (1 - r_h)$ представляет оставшуюся часть, обусловленную продуктами, не входящими в ассортимент. Наконец, вычислим спрос на один продукт как долю от общего спроса, пропорционального расчетному спросу на каждый продукт:

$$d_{jh} = \begin{cases} T_h \cdot r_h \cdot \frac{Q_{jh}}{\sum_{j \in N_h} Q_{jh}}, & \text{если } j \in N_h \\ T_h \cdot (1 - r_h) \cdot \frac{Q_{jh}^0}{\sum_{j \in N \setminus N_h} Q_{jh}^0}, & \text{если } j \in N \setminus N_h \end{cases}. \quad (6.135)$$

Все коэффициенты в уравнениях 6.135 и 6.132 можно получить по данным, поэтому все формулы можно свести к исходной задаче оптимизации 6.124 и решить ее с помощью численных методов.

Уравнение 6.124 дает множество предположительно оптимальных уровней запасов f_j для всех продуктов. Эти уровни можно использовать для корректировки запасов и оптимизации раскладки продуктов на полке. Важно отметить, что модель позволяет ретейлеру проводить анализ «что, если» для оценки влияния изменения ассортимента и уровня запасов на валовую прибыль. В частности, ретейлер может построить кривые, показывающие ожидаемую валовую прибыль, в зависимости

от уровня запасов для одного или группы продуктов. Такие кривые особенно описательны для скоропортящихся продуктов, поскольку валовая прибыль является выпуклой функцией, которая равна нулю, когда уровень запасов равен нулю, а также когда уровень запасов слишком высок, что влечет потери от утилизации просроченных продуктов, с максимумом между этими двумя крайностями.

6.10. Архитектура систем управления ценами

Разработка и внедрение алгоритмических систем управления ценами может существенно различаться в разных сферах, однако почти всегда управление ценами включает несколько основных процессов, которые можно рассматривать как функциональные компоненты эталонной логической архитектуры. На рис. 6.28 представлена обзорная диаграмма, где показаны эти ключевые компоненты и их взаимосвязи. Она включает три основные подсистемы.

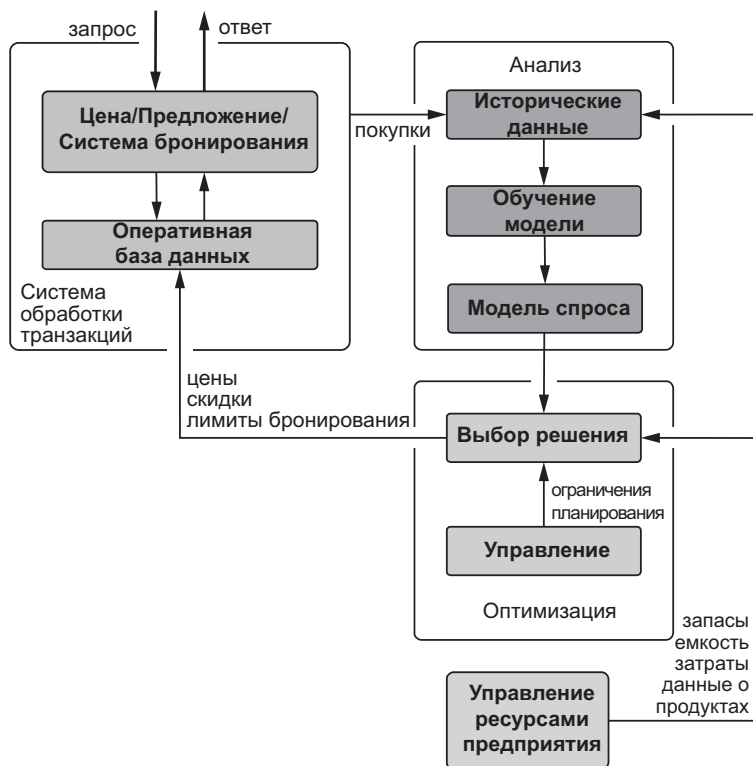


Рис. 6.28. Обобщенная архитектура системы управления ценами

СИСТЕМА ОБРАБОТКИ ТРАНЗАКЦИЙ. Цель системы обработки транзакций — исполнение ценовых решений, полученных от системы оптимизации. В некоторых средах оптимизированные цены и правила ценообразования могут загружаться в несколько систем обработки транзакций, которые вычисляют конечные цены независимо друг от друга. Например, ретейлер может обновлять цены в точках продаж (Point of Sales, POS) каждый день или неделю, а также обновлять цены в своей платформе электронной коммерции. Каждый магазин или платформа электронной коммерции работает независимо. В других средах может использоваться единая система обработки транзакций, централизованная или распределенная в техническом смысле, которая в режиме реального времени обрабатывает запросы на цены или квоты ресурсов. Такой подход, как правило, предпочтительнее, потому что сложные ценовые решения могут приниматься согласованным образом. Например, ретейлеры часто рассчитывают конечную цену для каждой сделки на основе прејскурантных цен, действующих скидок, номеров карт лояльности клиентов, введенных или отсканированных промокодов и других факторов. Системы бронирования в сфере услуг, например в авиакомпаниях или гостиницах, еще более интерактивны, потому что не только применяют правила ценообразования, но также отслеживают забронированные ресурсы и обновляют записи в оперативной базе данных после каждой транзакции.

Система обработки транзакций может принадлежать продавцу или быть общим ресурсом, предоставленным третьей стороной. Например, авиабилеты обычно бронируются через глобальные системы распределения (Global Distribution Systems, GDS), управляемые специализированными компаниями, и авиакомпании устанавливают лимиты бронирования непосредственно в GDS.

АНАЛИЗ. Прогнозирование спроса играет важную роль в управлении ценами, позволяя количественно оценивать цены и принимать решения по ассортименту. Это требует создания аналитической инфраструктуры, собирающей и хранящей исторические данные и поддерживающей обучение и применение предиктивных моделей. Мы видели, что модели прогнозирования спроса могут использовать довольно широкий спектр информации, включая данные о продуктах, скидках, конкурентных ценах, погодных данных и календарей праздничных дней. Многие из этих элементов данных можно извлекать из корпоративных систем управления ресурсами или получать от сторонних поставщиков данных.

Конвейер обучения модели спроса часто автоматизирован и работает в цикле или с определенным приращением, обновляя модель по мере поступления новых данных. Это гарантирует, что модели будут учитывать изменения в бизнесе и ландшафте конкуренции.

СИСТЕМА ОПТИМИЗАЦИИ. Основой системы оптимизации является компонент, выбирающий решения, который обычно использует численные методы

оптимизации для поиска ценового графика или уровней бронирования, максимизирующих доход. Компонент выбора решения настраивается административным компонентом управления, принимая от него параметры оптимизации и ограничения. В некоторых сферах, например в авиакомпаниях, процессом оптимизации может управлять команда аналитиков, которые следят за деятельностью системы и могут корректировать решения в особых случаях, например в преддверии крупных публичных мероприятий.

Подобно моделированию спроса, оптимизация также выполняется в цикле, пересчитывая цены или лимиты бронирования по мере обновления данных о продажах и запасах. Однако оптимизация, как правило, выполняется чаще, чем анализ, чтобы поспевать за быстро снижающимися уровнями емкости. Такая непрерывная оптимизация почти всегда используется для динамического ценообразования и распределения ресурсов, так же как другие методы оптимизации цен, которые не способны моделировать временные изменения напрямую, но постоянно корректирующие цены при многократном применении и, таким образом, становящиеся динамическими.

6.11. Итоги

- Ценовые решения чрезвычайно важны для конкурентоспособности и прибыльности фирмы. Совершенствование ценообразования часто оказывает гораздо большее влияние на прибыль, чем сопоставимая оптимизация продаж (каналы рекламы и сбыта), переменных затрат (поставки и производство) или постоянных затрат (операции и управление ресурсами).
- Цена является денежным выражением стоимости. Границы хорошей цены можно определить путем оценки потребительской полезности и сопоставимых альтернатив. Однако воспринимаемая ценность продукта или услуги зависит от того, как именно передается информация о продукте и цене. Например, эффективными мерами, как известно, являются снижение цен и разукрупнение пакетов. Это обеспечивает фундаментальное обоснование уценок и скидок.
- Спрос определяется готовностью клиентов платить. Различные распределения готовности платить дают различные кривые спроса. К основным кривым спроса относятся: линейная функция, функция постоянной эластичности и логит-функция.
- К основным ценовым структурам относятся: цена за единицу, сегментированная цена, двухкомпонентные тарифы, связывающие соглашения и пакетирование. Все эти структуры можно оптимизировать, если точно оценить кривую глобального спроса.

- Основные кривые спроса не учитывают сезонность, конкуренцию и свойства продукта. Алгоритмическая оптимизация цен требует разработки более мощных моделей спроса. Значительное число таких моделей можно найти в литературе по различным отраслям.
- Оптимизация цен обычно требует решения задачи численной оптимизации, адаптированной к конкретным структурным ограничениям, ограничениям спроса и предложения. Двумя основными бизнес-моделями оптимизации цен являются дифференциация цен, направленная на оптимизацию цен для нескольких сегментов рынка, и динамическое ценообразование, которое тоже можно рассматривать как метод сегментации рынка. К основным ограничениям, связанным с оптимизацией цен, относятся взаимозависимые функции спроса, ограниченная емкость и скоропортящиеся запасы.
- В сфере услуг, например в авиакомпаниях, гостиничном бизнесе, грузоперевозках и спорте, доли общей емкости могут резервироваться для разных классов тарифов. Такой подход к распределению ресурсов можно рассматривать как альтернативу динамическому ценообразованию. Наиболее простые методы распределения ресурсов позволяют оптимизировать уровни бронирования для отдельных единиц ограниченной емкости, таких как места в самолете, но существует большое количество более продвинутых методов, способных оптимизировать сети ресурсов и обрабатывать отмену бронирования.
- Оптимизация ассортимента тесно связана с управлением ценами, потому что обе задачи основаны на прогнозировании спроса. Оптимизация ассортимента нацелена на моделирование зависимостей между спросом на разные продукты и категории, что позволяет анализировать возможные изменения в ассортименте и перераспределении ресурсов, способные повлиять на прибыль.
- Цену, лимиты бронирования и ассортимент можно рассматривать как разные элементы управления, которые продавец может использовать для принятия решений о ценообразовании. Все методы оптимизации, связанные с этими элементами управления, используют в качестве базового строительного блока прогнозирования спроса, но выполняют разные бизнес-действия для его реализации.

Приложение.

Распределение Дирихле

Распределение Дирихле — относительно сложная тема, особенно в контексте применения в маркетинге, поэтому я решил дать краткое введение. Распределение Дирихле используется в главе 3, где рассматриваются неэкспериментальные исследования, а также в главе 4, где распределение Дирихле используется в тематическом моделировании.

В большинстве маркетинговых приложений приходится иметь дело с вероятностными распределениями случайных величин. Несмотря на то что распределение Дирихле можно рассматривать с этой точки зрения, оно имеет также существенно другое значение, что важно для наших целей. Рассмотрим упрощенный пример, иллюстрирующий этот аспект [Frigyik et al., 2010]. Шестигранный кубик можно рассматривать как дискретное распределение вероятности, генерирующее числа от одного до шести. С совершенным кубиком все числа имеют одинаковую вероятность выпадения, равную одной шестой. Однако распределение вероятности для реального кубика будет отличаться от равномерного из-за несовершенства производства и других физических факторов. Если взять мешок со 100 кубиками, каждому кубику в нем будет соответствовать своя *функция распределения масс* (Probability Mass Function, PMF), а мешок кубиков соответствует распределению функций PMF. Свойства этого распределения зависят от качества изготовления: распределение функций PMF может существенно отклоняться от равномерного в случае низкогокачественного производства, или они могут быть почти идентичны, если выдерживается высокая точность при изготовлении. Это распределение функций PMF можно описать с помощью распределения Дирихле.

Более практичным и актуальным для нас примером может служить коллекция текстовых документов. Учитывая, что в сумме документы содержат m разных слов,

каждый документ можно рассматривать как функцию РМФ, которую можно оценить путем подсчета частоты появления каждого слова в документах. Коллекция документов — это коллекция функций РМФ, и мы можем подобрать параметры распределения Дирихле для этой коллекции. Выражаясь формальным языком, каждый документ d можно смоделировать как вектор вероятностей m слов, сумма которых должна быть равна единице:

$$p_{d1} + \dots + p_{dm} = 1, \quad p_{di} \in [0, 1]. \quad (\text{П.1})$$

Геометрически это уравнение описывает $(m - 1)$ -мерный симплекс в m -мерном пространстве. Например, коллекция документов с тремя разными словами соответствует двумерному треугольнику (симплексу) в трехмерном пространстве, как показано на рис. П.1. Каждая точка симплекса соответствует допустимой функции РМФ, тогда как все другие точки пространства не соответствуют никакой допустимой РМФ. Мы можем сгенерировать коллекцию из m документов, указав распределение по симплексу, выбрав m функций РМФ из этого распределения и сгенерировав m -й документ, выбирая термы из соответствующей РМФ.

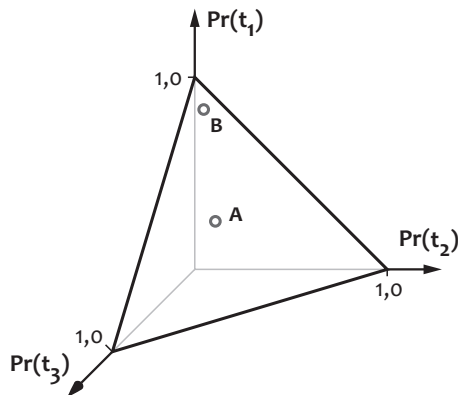


Рис. П.1. Трехмерный симплекс вероятности. Точка А соответствует документу с равномерным распределением термов; точка В соответствует документу, где терм t_1 гораздо более вероятен, чем два других терма

Распределение Дирихле является распределением вероятности по симплексу. Для m измерений каждый экземпляр, взятый из распределения Дирихле, является m -компонентной функцией распределения масс:

$$\mathbf{p} = (p_1, \dots, p_m), \quad \sum_i p_i = 1. \quad (\text{П.2})$$

Само распределение задается вектором m параметров:

$$\alpha = (\alpha_1, \dots, \alpha_m), \quad \alpha_i > 0, \quad (\text{П.3})$$

где каждый параметр можно рассматривать как вес соответствующего компонента. Функция плотности вероятности распределения Дирихле определяется как

$$\text{Dir}(\alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^m p_i^{\alpha_i-1}, \quad (\text{П.4})$$

где $B(\alpha)$ — константа нормализации, определяемая выражением

$$B(\alpha) = \frac{\prod_{i=1}^m \Gamma(\alpha_i)}{\Gamma\left(\sum_{i=1}^m \alpha_i\right)}. \quad (\text{П.5})$$

Если все элементы вектора параметров имеют одинаковое значение, распределение полностью определяется этим единственным значением, называемым *параметром концентрации*. На рис. П.2 изображены функции плотности вероятности для распределения Дирихле в трехмерном пространстве с разными значениями параметров.

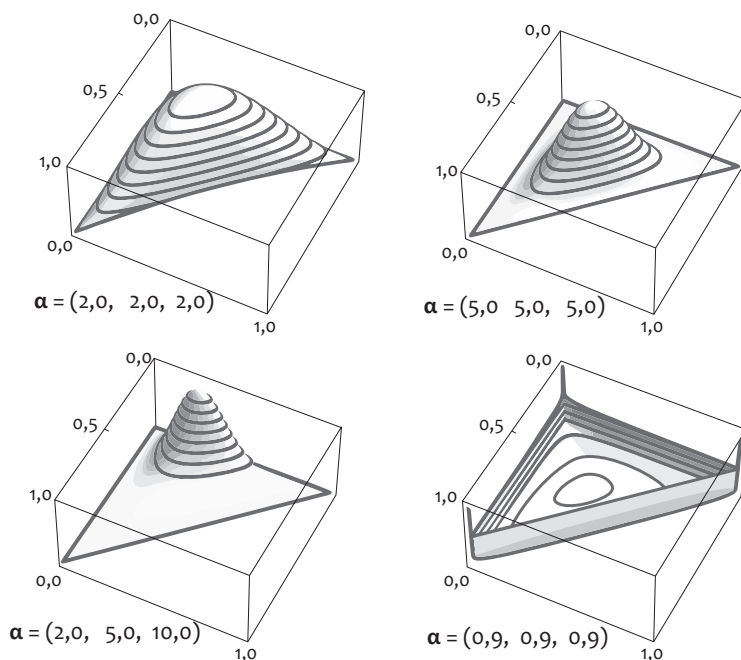


Рис. П.2. Графики плотности для распределения Дирихле по симплексу вероятности в трехмерном пространстве

Функция плотности имеет плоский вид для $\alpha = (1, 1, 1)$; колоколообразную форму, симметричную относительно центра симплекса, когда все элементы вектора параметров равны. Если элементы вектора не равны, колокол смещается в сторону параметров с наибольшими величинами. Наконец, важно отметить, что распределение Дирихле является разреженным для параметров с малыми величинами, в том смысле что плотность сосредоточена в углах, и, следовательно, функции РМФ, взятые из такого распределения, как правило, имеют сильный уклон в сторону небольшого подмножества термов [Telgarsky, 2013].

Библиография

Adomavicius, G. and Tuzhilin, A. (2008). Context-aware recommender systems. In *Proceedings of the 2008 ACM Conference on Recommender Systems*, RecSys '08, pages 335–336, New York, NY, USA. ACM.

Aggarwal, C. C. (2016). *Recommender Systems: The Textbook*. Springer Publishing Company, Incorporated.

Almon, S. (1965). The distributed lag between capital appropriations and expenditures. *Econometrica*, 33(1):178–196.

Anderson, C. (2008). *The Long Tail: Why the Future of Business is Selling Less of More*. Hyperion.¹

Anupindi, R., Dada, M., and Gupta, S. (1998). Estimation of consumer demand with stock-out based substitution: An application to vending machine products. *Marketing Science*, 17(4):406–423.

Artun, O. and Levin, D. (2015). *Predictive Marketing: Easy Ways Every Marketer Can Use Customer Analytics and Big Data*. Wiley.

Asuncion, A., Welling, M., Smyth, P., and Teh, Y. W. (2009). On smoothing and inference for topic models. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, UAI '09, pages 27–34. AUAI Press.

Bell, R. M. and Koren, Y. (2007). Scalable collaborative filtering with jointly derived neighborhood interpolation weights. In *Proceedings of the 2007 Seventh IEEE International Conference on Data Mining*, ICDM '07, pages 43–52. IEEE Computer Society.

Bellogin, A., Castells, P. and Cantador, I. (2014). Neighbor selection and weighting in user-based collaborative filtering: A performance prediction approach. *ACM Trans. Web*, 8(2):12:1–12:30,

Belobaba, P. (1987). Air travel demand and airline seat inventory management. Technical report, Cambridge, MA: Flight Transportation Laboratory, Massachusetts Institute of Technology.

Belobaba, P. P. (1989). Application of a probabilistic decision model to airline seat inventory control. *Operations Research*, 37(2):183–197.

¹ Крис Андерсон. *Длинный хвост. Новая модель ведения бизнеса*. М., Вершина, 2008. — Примеч. пер.

- Bergamaschi, S., Po, L. and Sorrentino, S. (2014). Comparing topic models for a movie recommendation system.
- Berger, P. and Nasr, N. (1998). *Customer Lifetime Value: Marketing Models and Applications*, volume 12.
- Berry, M. (2009). *Differential Response or Uplift Modeling*.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- Blattberg, R. C. and Deighton, J. A. (1996). *Manage Marketing by the Customer Equity Test*, volume 74.
- Blei, D. M., Ng, A. Y. and Jordan, M. I. (2003). Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022.
- Bradford, R. (2008). An empirical study of required dimensionality for large-scale latent semantic indexing applications. In *Proceedings of the 17th ACM Conference on Information and Knowledge Management*. ACM.
- Breese, J. S., Heckerman, D. and Kadie, C. (1998). *Empirical Analysis of Predictive Algorithms for Collaborative Filtering*. UAI'98. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- Breslow, N. E. (1972). Discussion following «regression models and life tables». *Journal of the Royal Statistical Society*, 34:187–220,
- Breslow, N. E. (1974). Covariance analysis of censored survival data. *Biometrics*, 30(1):89–99.
- Broadbent, S. (1979). One way tv advertisements work. *Journal of the Market Research Society*, 23(3).
- Brumelle, S. L. and McGill, J. I. (1993). *Airline Seat Allocation with Multiple Nested Fare Classes*, volume 41.
- Burges, C., Shaked, T., Renshaw, E., Lazier, A., Deeds, M., Hamilton, N. and Hullender, G. (2005). Learning to rank using gradient descent. In *Proceedings of the 22Nd International Conference on Machine Learning*, ICML '05.
- Burges, C. J. C. (2010). From RankNet to LambdaRank to LambdaMART: An overview. Technical report, Microsoft Research.
- Burke, R. (2002). Hybrid recommender systems: Survey and experiments. *User Modeling and User-Adapted Interaction*, 12(4):331–370,
- Cao, Z., Qin, T., Liu, T.-Y., Tsai, M.-F. and Li, H. (2007). Learning to rank: From pairwise approach to listwise approach. In *Proceedings of the 24th International Conference on Machine Learning*, pages 129–136. ACM.
- Caro, F. and Gallien, J. (2012). *Clearance Pricing Optimization for a Fast-Fashion Retailer*, volume 60.
- Carpenter, G. S. and Shankar, V. (2013). *Handbook of Marketing Strategy*. Elgar original reference. Elgar.
- Catalina Marketing (2014). *Catalina Category Marketing*.
- Chalasani, P. and Sriharsha, R. (2016). *Monte Carlo Simulations in Ad Lift Measurement Using Spark*.
- Chapelle, O. and Chang, Y. (2011). Yahoo! learning to rank challenge overview. *Journal of Machine Learning Research – Proceedings Track*, 14:1–24.
- Chickering, D. M. and Pearl, J. (1996). *A Clinician's Tool for Analyzing Non-compliance*.

- Cossock, D. and Zhang, T. (2006). Subset ranking using regression. In *Proceedings of the 19th Annual Conference on Learning Theory*, pages 605–619. Springer-Verlag.
- Cox, D. R. (1972). Regression models and life-tables. *Journal of the Royal Statistical Society*, 34(2):187–220,
- Cox, D. R. (1975). Partial likelihood. *Biometrika*, 62(2):269–276.
- Crammer, K. and Singer, Y. (2001). Pranking with ranking. In *Advances in Neural Information Processing Systems 14*, pages 641–647. MIT Press.
- Cremonesi, P., Koren, Y. and Turrin, R. (2010). Performance of recommender algorithms on top-n recommendation tasks. In *Proceedings of the Fourth ACM Conference on Recommender Systems*, pages 39–46. ACM.
- Crocker, C., Kulick, A. and Ram, B. (2012). Real user monitoring at walmart.
- Cuellar, S. S. and Brunamonti, M. (2013). *Retail Channel Price Discrimination*, volume 21.
- Dalessandro, B., Perlich, C., Hook, R., Stitelman, O., Raeder, T. and Provost, F. (2012a). *Bid Optimizing and Inventory Scoring in Targeted Online Advertising*.
- Dalessandro, B., Perlich, C., Stitelman, O. and Provost, F. (2012b). *Causally Motivated Attribution for Online Advertising*.
- Debreu, G. (1960). *Review of R. D. Luce, Individual Choice Behavior: A Theoretical Analysis*.
- Deerwester, S., Dumais, S. T., Furnas, G.W., Landauer, T. K., and Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6).
- Devooght, R., Kourtellis, N., and Mantrach, A. (2015). Dynamic matrix factorization with priors on unknown values. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 189–198. ACM.
- Ding, Y. and Li, X. (2005). Time weight collaborative filtering. In *Proceedings of the 14th ACM International Conference on Information and Knowledge Management*, CIKM '05, pages 485–492, New York, NY, USA. ACM.
- Duhigg, C. (2012). *How Companies Learn Your Secrets*.
- Efron, B. (1977). The efficiency of cox's likelihood function for censored data. *Journal of the American Statistical Association*, 72(359):557–565.
- Ekstrand, M. D., Riedl, J. T., and Konstan, J. A. (2011). Collaborative filtering recommender systems. *Found. Trends Hum.-Comput. Interact.*, 4(2):81–173.
- Falk, K. (2017). *Practical Recommender Systems*. Manning Publications Company.
- Ferreira, K. J., Lee, B. H. A. and Simchi-Levi, D. (2016). *Analytics for an Online Retailer: Demand Forecasting and Price Optimization*, volume 18.
- Freund, Y., Iyer, R., Schapire, R. E. and Singer, Y. (2003). An efficient boosting algorithm for combining preferences. *J. Mach. Learn. Res.*, 4:933–969.
- Frigyik, B., Kapila, A., and Maya, G. (2010). Introduction to the dirichlet distribution and related processes. Technical report, University of Washington.
- Funk, S. (2016). Netflix update: Try this at home.
- Ge, M., Delgado-Battenfeld, C. and Jannach, D. (2010). Beyond accuracy: Evaluating recommender systems by coverage and serendipity. In *Proceedings of the Fourth ACM Conference on Recommender Systems*, pages 257–260, ACM.

- Geman, S. and Geman, D. (1984). *Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images*, volume 6. IEEE Computer Society, Washington, DC, USA.
- Geng, L. and Hamilton, H. J. (2006). Interestingness measures for data mining: A survey. *ACM Comput. Surv.*, 38(3).
- Ghani, R. and Fano, A. (2002). Building recommender systems using a knowledge base of product semantics. In *2nd International Conference on Adaptive Hypermedia and Adaptive Web Based Systems, Malaga*.
- Ghani, R., Probst, K., Liu, Y., Krema, M. and Fano, A. (2006). Text mining for product attribute extraction. *SIGKDD Explor. Newsl.*, pages 41–48.
- Giunchiglia, F., Kharkevich, U. and Zaihrayeu, I. (2009). Concept search. In *Proceedings of the 6th European Semantic Web Conference on The Semantic Web: Research and Applications, ESWC 2009 Heraklion*. Springer-Verlag.
- Goldberg, D., Nichols, D., Oki, B. M. and Terry, D. (1992). Using collaborative filtering to weave an information tapestry. *Commun. ACM*, 35(12):61–70.
- Google Inc. (2011). *The Arrival of Real-Time Bidding*.
- Gormley, C. and Tong, Z. (2015). *Elasticsearch: The Definitive Guide*. O'Reilly Media.
- Green, P. E. and Srinivasan, V. (1978). *Conjoint Analysis in Consumer Research: Issues and Outlook*, volume 5.
- Grigsby, M. (2016). *Advanced Customer Analytics: Targeting, Valuing, Segmenting and Loyalty Techniques*. Marketing Science Series. Kogan Page, Limited.
- Hall, R. E. (1967). Polynomial distributed lags.
- Heckerman, D. and Shachter, R. D. (1995). *Decision-Theoretic Foundations for Causal Reasoning*.
- Herbrich, R., Graepel, T. and Obermayer, K. (2000). Large margin rank boundaries for ordinal regression. In *Advances in Large Margin Classifiers*, pages 115–132. MIT Press.
- Herlocker, J. L., Konstan, J. A., Borchers, A. and Riedl, J. (1999). An algorithmic framework for performing collaborative filtering. In *Proceedings of the 22Nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '99*, pages 230–237, New York, NY, USA. ACM.
- Herlocker, J. L., Konstan, J. A., Terveen, L. G. and Riedl, J. T. (2004). Evaluating collaborative filtering recommender systems. pages 5–53.
- Hofmann, T. (1999). Probabilistic latent semantic indexing. In *Proceedings of the 22Nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '99*, pages 50–57, New York, NY, USA. ACM.
- Hu, Y., Koren, Y. and Volinsky, C. (2008). Collaborative filtering for implicit feedback datasets. In *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining, ICDM '08*, pages 263–272. IEEE Computer Society.
- Hughes, S. (2015). Implementing conceptual search in Solr using LSA and Word2Vec.
- Jack, K., Ingold, E. and Hristakeva, M. (2016). Mendeley suggest architecture.
- Jahrer, M., Töschner, A. and Legenstein, R. (2010). Combining predictions for accurate recommender systems. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '10*, pages 693–702, New York, NY, USA. ACM.

- Jambor, T. and Wang, J. (2010). Optimizing multiple objectives in collaborative filtering. In *Proceedings of the Fourth ACM Conference on Recommender Systems, RecSys '10*, pages 55–62, New York, NY, USA. ACM.
- Järvelin, K. and Kekäläinen, J. (2000). IR evaluation methods for retrieving highly relevant documents. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*.
- Jo, B. (2002). *Statistical Power in Randomized Intervention Studies With Noncompliance*.
- Johnson, J., Tellis, G. J. and Ip, E. H. (2013). *To Whom, When, and How Much to Discount? A Constrained Optimization of Customized Temporal Discounts*, volume 89. Elsevier.
- Ju, C., Bao, F., Xu, C., and Fu, X. (2015). A novel method of interestingness measures for association rules mining based on profit. In *Discrete Dynamics in Nature and Society*.
- Kahneman, D. and Tversky, A. (1979). *Prospect theory: An analysis of decisions under risk*.
- Kamotsky, D. and Vargas, M. (2014). System and method for performing a pattern matching search. US Patent App. 14/292,018.
- Kane, K., Lo, V. S. and Zheng, J. X. (2014). Mining for the truly responsive customers and prospects using true-lift modeling: Comparison of new and existing methods. *Journal of Marketing Analytics*, 2(4):218–238.
- Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, 53(282):457–481.
- Khan, R. J. and Jain, D. C. (2005). *An Empirical Analysis of Price Discrimination Mechanisms and Retailer Profitability*, volume 42.
- Khludnev, M. (2013). *Concept search for eCommerce with Solr*.
- Kleinberg, J., Papadimitriou, C. and Raghavan, P. (1998). *A Microeconomic View of Data Mining*, volume 2. Kluwer Academic Publishers, Hingham, MA, USA.
- Kohavi, R. and Longbotham, R. (2007). Online experiments: Lessons learned.
- Kök, A. G. and Fisher, M. L. (2007). *Demand Estimation and Assortment Optimization Under Substitution: Methodology and Application*, volume 55. INFORMS.
- Koren, Y. (2007). How useful is a lower rmse?
- Koren, Y. (2008). Factorization meets the neighborhood: A multifaceted collaborative filtering model. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 426–434, New York, NY, USA. ACM.
- Koren, Y. (2009). The bellkor solution to the netflix grand prize.
- Koren, Y. and Bell, R. M. (2011). Advances in collaborative filtering. In Ricci, F., Rokach, L., Shapira, B. and Kantor, P. B., editors. *Recommender Systems Handbook*, pages 145–186. Springer.
- Koyck, L. M. (1954). Distributed lags and investment analysis.
- Lauterborn, B. (1990). *New Marketing Litany: Four Ps Passe: C-Words Take Over*.
- Lee, D. D. and Seung, H. S. (2001). Algorithms for non-negative matrix factorization. In Leen, T. K., Dietterich, T. G., and Tresp, V., editors, *Advances in Neural Information Processing Systems 13*, pages 556–562. MIT Press.
- Li, P., Burges, C. J. C. and Wu, Q. (2007). Mcrank: Learning to rank using multiple classification and gradient boosting. In *NIPS*, pages 897–904. Curran Associates, Inc.

- Linden, G., Smith, B. and York, J. (2003). Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Computing*, 7(1):76–80.
- Littlewood, K. (1972). *Forecasting and Control of Passenger Bookings*.
- Liu, T.-Y. (2009). Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval*, (3):225–331.
- Liu, T.-Y. and Qin, T. (2010). Microsoft learning to rank datasets.
- Lo, V. S. (2002). The true lift model: a novel data mining approach to response modeling in database marketing. *ACM SIGKDD Explorations Newsletter*, 4(2):78–86.
- Manning, C. D., Raghavan, P. and Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA.¹
- Marn, M. and Roseillo, R. (1992). *Managing price, gaining profit*.
- McCarthy, E. J. (1960). *Basic Marketing. A Managerial Approach*.
- Melville, P., Mooney, R. J. and Nagarajan, R. (2002). Content-boosted collaborative filtering for improved recommendations. In *Eighteenth National Conference on Artificial Intelligence*, pages 187–192. American Association for Artificial Intelligence.
- Mikolov, T., Chen, K., Corrado, G. and Dean, J. (2013a). Efficient estimation of word representations in vector space.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S. and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems 26*. Curran Associates, Inc.
- Miyahara, K. and Pazzani, M. J. (2000). Collaborative filtering with the simple bayesian classifier. In *Proceedings of the 6th Pacific Rim International Conference on Artificial Intelligence, PRICAI'00*, pages 679–689. Springer-Verlag.
- Mobasher, B., Dai, H., Luo, T. and Nakagawa, M. (2001). Effective personalization based on association rule discovery from web usage data. In *Proceedings of the 3rd International Workshop on Web Information and Data Management, WIDM '01*, pages 9–15, New York, NY, USA. ACM.
- Mooney, R. J. and Roy, L. (1999). Content-based book recommending using learning for text categorization. In *Proceedings of the SIGIR-99 Workshop on Recommender Systems: Algorithms and Evaluation*.
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. The MIT Press.
- Musalem, A., Olivares, M., Bradlow, E. T., Terwiesch, C. and Corsten, D. (2010). Structural estimation of the effect of out-of-stocks. *Management Science*, 56(7):1180–1197.
- Oi, W. Y. (1971). *A Disneyland Dilemma: Two-Part Tariffs for a Mickey Mouse Monopoly*.
- Oneata, D. (1999). Probabilistic latent semantic analysis. In *Proceedings of the Fifteenth conference on Uncertainty*, pages 1–7.
- Park, L. A. and Ramamohanarao, K. (2009). Efficient storage and retrieval of probabilistic latent semantic information for information retrieval. *The VLDB Journal*, 18(1):141–155.
- Pashigian, P. (1987). *Demand Uncertainty and Sales: A Study of Fashion and Markdown Pricing*.

¹ Кристофер Д. Маннинг, Прабхакар Рагхаван, Хайнрих Шютце. Введение в информационный поиск. М., Вильямс, 2014. — *Примеч. пер.*

- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *Conference on Empirical Methods on Natural Language Processing*, volume 14, pages 1532–1543.
- Perlich, C., Dalessandro, B., Raeder, T., Stitelman, O., and Provost, F. (2013). *Machine Learning for Targeted Display Advertising: Transfer Learning in Action*.
- Peter, G. and Eugene, S. (2015). Deep data at macys: Searching hierarchical documents for ecommerce merchandising.
- Pfeifer, P. and Carraway, R. (2000). *Modeling Customer Relationships as Markov Chains*.
- Porter, M. F. (1980). An algorithm for suffix stripping. *Program*, 14(3):130–137.
- Radcliffe, N. J. and Simpson, R. (2007). Identifying who can be saved and who will be driven away by retention activity.
- Radcliffe, N. J. and Surry, P. (1999). Differential response analysis: Modeling true response by isolating the effect of a single action. *Credit Scoring and Credit Control VI. Edinburgh, Scotland*.
- Radcliffe, N. J. and Surry, P. D. (2011). Real-world uplift modeling with significance-based uplift trees. *White Paper TR-2011-1, Stochastic Solutions*.
- Radlinski, F. and Joachims, T. (2005). Query chains: Learning to rank from implicit feedback. In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, pages 239–248. ACM.
- Rodriguez, M., Posse, C. and Zhang, E. (2012). Multiple objective optimization in recommender systems. In *Proceedings of the Sixth ACM Conference on Recommender Systems, RecSys '12*, pages 11–18, New York, NY, USA. ACM.
- Rohde, D. L. T., Gonnerman, L. M. and Plaut, D. C. (2006). An improved model of semantic similarity based on lexical co-occurrence. *Communications of the ACM*.
- Rong, X. (2014a). wevi: Word embedding visual inspector.
- Rong, X. (2014b). Word2Vec parameter learning explained.
- Rubin, D. (1974). *Estimating causal effects of treatments in randomized and nonrandomized studies*.
- Sarwar, B., Karypis, G., Konstan, J. and Riedl, J. (2001). Item-based collaborative filtering recommendation algorithms. In *Proceedings of the 10th International Conference on World Wide Web*, pages 285–295. ACM.
- Shao, X. and Li, L. (2011). *Data-driven Multi-touch Attribution Models*.
- Sill, J., Takacs, G., Mackey, L. W. and Lin, D. (2009). Feature-weighted linear stacking.
- Smith, B., Leimkuhler, J. and Darrow, R. (1992). *Yield Management in American Airlines*.
- Smith, T. (2012). *Pricing Strategy*. South-Western Cengage Learning.
- Spangher, A. (2015). Building the next new york times recommendation engine.
- Su, X. and Khoshgoftaar, T. M. (2006). Collaborative filtering for multiclass data using belief nets algorithms. In *Proceedings of the 18th IEEE International Conference on Tools with Artificial Intelligence, ICTAI '06*, pages 497–504, Washington, DC, USA. IEEE Computer Society.
- Su, X. and Khoshgoftaar, T. M. (2009). A survey of collaborative filtering techniques. *Adv. in Artif. Intell.*, pages 4:2–4:2.

- Su, X., Khoshgoftaar, T. M., Zhu, X. and Greiner, R. (2008). Imputation-boosted collaborative filtering using machine learning classifiers. In *Proceedings of the 2008 ACM Symposium on Applied Computing*, SAC '08, pages 949–950, New York, NY, USA. ACM.
- Talluri, K. and Van Ryzin, G. (2004). *The Theory and Practice of Revenue Management*. International Series in Operations Research & Management Science. Springer.
- Telgarsky, M. (2013). Dirichlet draws are sparse with high probability.
- Terry, D. B. (1993). A tour through tapestry. In *Proceedings of the Conference on Organizational Computing Systems*, COCS '93, pages 21–30, New York, NY, USA. ACM.
- Töscher, A., Jahrer, M. and Bell, R. M. (2009). The bigchaos solution to the netflix grand prize.
- Train, K. (2003). *Discrete Choice Methods with Simulation*. SUNY-Oswego, Department of Economics.
- Turnbull, D. and Berryman, J. (2016). *Relevant Search. With applications for Solr and Elasticsearch*. Manning Publications.¹
- Vasigh, B., Tacker, T. and Fleming, M. (2013). *Introduction to Air Transport Economics: From Theory to Applications*.
- Vulcano, G. J., van Ryzin, G. J. and Ratliff, R. (2012). Estimating primary demand for substitutable products from sales transaction data. *Operations Research*, 60(2):313–334.
- Walker, R. (2009). *The Song Decoders*. The New York Times Magazine.
- Wierenga, B. (2010). The interface of marketing and operations research. In Kroon, L., Zuidwijk, R., and Li, T., editors, *Liber Amicorum in Memoriam Jo van Nunen*.
- Winston, W. L. (2014). *Marketing Analytics: Data-Driven Techniques with Microsoft Excel*. Wiley Publishing, 1st edition.
- Xia, Z., Dong, Y., and Xing, G. (2006). Support vector machines for collaborative filtering. In *Proceedings of the 44th Annual Southeast Regional Conference*, ACM-SE 44, pages 169–174, New York, NY, USA. ACM.
- Xu, J. and Li, H. (2007). Adarank: A boosting algorithm for information retrieval. In *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 391–398, New York, NY, USA. ACM.
- Zaki, M. J. and Meira, W. (2014). *Data Mining and Analysis: Fundamental Concepts and Algorithms*.
- Zhang, A., Goyal, A., Kong, W., Deng, H., Dong, A., Chang, Y., Gunter, C. A. and Han, J. (2015). adaqac: Adaptive query auto-completion via implicit negative feedback. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 143–152.
- Zhang, H. (2004). The optimality of naive bayes. In Barr, V. and Markov, Z., editors. *Proceedings of the Seventeenth International Florida Artificial Intelligence Research Society Conference (FLAIRS 2004)*. AAAI Press.
- Zhang, S., Wang, W., Ford, J. and Makedon, F. (1996). Learning from incomplete ratings using non-negative matrix factorization. In *In Proc. of the 6th SIAM Conference on Data Mining*, pages 549–553.

¹ Джон Берримен, Даг Тарнбулл. Релевантный поиск с использованием Elasticsearch и Solr. М., ДМК-Пресс. — *Примеч. пер.*

Илья Кацов

Машинное обучение для бизнеса и маркетинга

Перевел с английского А. Киселев

Заведующая редакцией	<i>Ю. Сергиенко</i>
Ведущий редактор	<i>К. Тульцева</i>
Литературный редактор	<i>А. Бульченко</i>
Художественный редактор	<i>В. Мостипан</i>
Корректоры	<i>Н. Сидорова, Г. Шкатова</i>
Верстка	<i>Л. Егорова</i>

Изготовлено в России. Изготовитель: ООО «Прогресс книга».
Место нахождения и фактический адрес: 194044, Россия, г. Санкт-Петербург,
Б. Сапсониевский пр., д. 29А, пом. 52. Тел.: +78127037373.

Дата изготовления: 06.2019. Наименование: книжная продукция. Срок годности: не ограничен.

Налоговая льгота — общероссийский классификатор продукции ОК 034-2014, 58.11.12 — Книги печатные профессиональные, технические и научные.

Импортер в Беларусь: ООО «ПИТЕР М», 220020, РБ, г. Минск, ул. Тимирязева, д. 121/3, к. 214, тел./факс: 208 80 01.

Подписано в печать 29.05.19. Формат 70×100/16. Бумага офсетная. Усл. п. л. 41,280. Тираж 1500. Заказ 0000.

Отпечатано в полном соответствии с качеством предоставленных издательством материалов
в ОГУП «Областная типография «Печатный двор». 432049, Ульяновск, Пушкирева д.27.