

МАШИННОЕ ОБУЧЕНИЕ НА R

 ЭКСПЕРТНЫЕ ТЕХНИКИ
ДЛЯ ПРОГНОСТИЧЕСКОГО АНАЛИЗА

 Packt

Бретт Ланц

Машинное обучение на R: экспертные техники для
прогностического анализа

Питер
2020

Об авторе

Бретт Ланц (Brett Lantz, @DataSpelunking) более десяти лет использует инновационные методы обработки данных для изучения поведения человека. Будучи по образованию социологом, Бретт впервые увлекся машинным обучением во время исследования большой базы профилей подростков в социальных сетях. Бретт — преподаватель DataCamp и часто выступает с докладами на конференциях и семинарах по машинному обучению по всему миру. Он известный энтузиаст в сфере практического применения науки о данных в области спорта, беспилотных транспортных средств, изучения иностранных языков и моды, а также во многих других отраслях. Бретт надеется в один прекрасный день написать обо всем этом на сайте dataspelunking.com, посвященном обмену знаниями о поиске закономерностей в данных.

Я не смог бы написать эту книгу без поддержки моей семьи. В частности, моя жена Джессика заслуживает огромной благодарности за ее бесконечное терпение и поддержку. Мои сыновья Уилл и Кэл родились в тот период, когда создавались первое и второе издания соответственно, и я бы не смог написать третье, если бы они меня отвлекали. Я посвящаю им эту книгу в надежде, что однажды она вдохновит их на решение больших задач. Желаю им следовать своему любопытству, куда бы оно их ни привело.

Я также признателен многим другим людям, которые косвенно поддержали эту книгу. Общение с педагогами, коллегами и сотрудниками Мичиганского университета, Университета Нотр-Дам и Университета Центральной Флориды способствовало рождению многих идей, которые я попытался выразить в тексте; что же касается отсутствия ясности в их выражении, то это сугубо мое упущение. Кроме того, эта книга могла бы вообще не появиться без более широкого сообщества исследователей, которые поделились своим опытом в виде статей, лекций и исходного кода. Наконец, я ценю усилия команд R и RStudio, а также всех тех, кто внес вклад в создание R-пакетов. Благодаря проделанной работе мы смогли донести идеи машинного обучения до широкой публики. Я искренне надеюсь, что мой труд также станет важной частью этой мозаики.

О научном редакторе

Рагхав Бали (Raghav Bali) — старший научный сотрудник одной из крупнейших в мире организаций здравоохранения. Занимается исследованиями и разработкой корпоративных решений, основанных на машинном обучении, глубоком обучении и обработке естественного языка для использования в области здравоохранения и страхования. На своей предыдущей должности в Intel он участвовал в реализации проактивных

инициатив в области информационных технологий, основанных на больших данных, с использованием обработки естественного языка, глубокого обучения и традиционных статистических методов. В American Express работал в области цифрового взаимодействия и удержания клиентов.

Рагхав является автором нескольких книг, выпущенных ведущими издательствами. Его последняя книга посвящена новейшим достижениям в области исследования трансферного обучения.

Рагхав окончил Международный институт информационных технологий в Бангалоре, имеет степень магистра (диплом с отличием). В те редкие моменты, когда он не занят решением научных проблем, Рагхав любит читать и фотографировать все подряд.

Предисловие

В основе машинного обучения (англ. Machine Learning, ML) лежат алгоритмы, которые преобразуют информацию в практически ценные данные. Именно поэтому машинное обучение так популярно в современную эру больших данных. Без него было бы почти невозможно отслеживать огромный поток информации.

Учитывая растущую популярность R — кросс-платформенной статистической свободно распространяемой среды программирования, — еще никогда не было более подходящего времени, чтобы начать использовать машинное обучение. R предоставляет мощный, но простой в освоении набор инструментов, которые помогут вам постигнуть суть ваших данных.

Сочетая практические примеры с базовой теорией, которая требуется для понимания того, как все работает внутри, эта книга даст вам возможность получить все необходимые знания, чтобы можно было начать работу с машинным обучением.

Для кого предназначена книга

Книга предназначена для тех, кто рассчитывает использовать данные в конкретной области. Возможно, вы уже немного знакомы с машинным обучением, но никогда не работали с языком R; или, наоборот, немного знаете об R, но почти не знаете о машинном обучении. В любом случае эта книга поможет вам быстро начать работу. Было бы полезно немного освежить в памяти основные понятия математики и программирования, но никакого предварительного опыта не потребуется. Вам нужно лишь желание учиться.

О чем вы прочтете в издании

Глава 1 «Введение в машинное обучение» содержит терминологию и понятия, которые определяют и выделяют теорию машинного обучения среди других областей, а также включает информацию о том, как выбрать алгоритм, подходящий для решения конкретной задачи.

Глава 2 «Управление данными и их интерпретация» даст вам возможность полностью погрузиться в работу с данными в среде R. Здесь речь пойдет об основных структурах данных и процедурах, используемых для загрузки, исследования и интерпретации данных.

Глава 3 «Ленивое обучение: классификация с использованием метода ближайших соседей» научит вас понимать и применять простой, но мощный алгоритм машинного обучения для решения вашей первой практической задачи: выявления особо опасных видов рака.

Глава 4 «Вероятностное обучение: классификация с использованием наивного байесовского классификатора» раскрывает основные понятия теории вероятностей, которые используются в современных системах фильтрации спама. Создавая собственный фильтр спама, вы изучите основы интеллектуального анализа текста.

Глава 5 «Разделяй и властвуй: классификация с использованием деревьев решений и правил» посвящена нескольким обучающим алгоритмам, прогнозы которых не только точны, но и легко интерпретируемы. Мы применим эти методы к задачам, в которых важна прозрачность.

Глава 6 «Прогнозирование числовых данных: регрессионные методы» познакомит с алгоритмами машинного обучения, используемыми для числовых прогнозов. Поскольку эти методы тесно связаны с областью статистики, вы также изучите базовые понятия, необходимые для понимания числовых отношений.

Глава 7 «Методы “черного ящика”: нейронные сети и метод опорных векторов» описывает два сложных, но мощных алгоритма машинного обучения. Их математика на первый взгляд может вас испугать, однако мы разберем примеры, иллюстрирующие их внутреннюю работу.

Глава 8 «Обнаружение закономерностей: анализ потребительской корзины с помощью ассоциативных правил» объясняет алгоритм, используемый в рекомендательных системах, применяемых во многих компаниях розничной торговли. Если вы когда-нибудь задумывались о том, почему системы розничных продаж знают ваши покупательские привычки лучше, чем вы сами, то эта глава раскроет их секреты.

Глава 9 «Поиск групп данных: кластеризация методом k-средних» посвящена процедуре поиска кластеров связанных элементов. Мы воспользуемся этим алгоритмом для идентификации профилей в онлайн-сообществе.

Глава 10 «Оценка эффективности модели» предоставит информацию о том, как измерить успешность проекта машинного обучения и получить надежный прогноз использования конкретного метода в будущем на других данных.

Глава 11 «Повышение эффективности модели» раскрывает методы, используемые теми, кто возглавляет список лидеров в области машинного обучения. Если в вас живет дух соревновательности или вы просто хотите получить максимальную отдачу от своих данных, то вам необходимо добавить эти методы в свой арсенал.

Глава 12 «Специальные разделы машинного обучения» исследует границы машинного обучения: от обработки больших данных до ускорения работы R. Прочитав ее, вы откроете для себя новые горизонты и узнаете, что еще можно делать с помощью R.

Что вам нужно для чтения книги

Примеры в этой книге написаны и протестированы для версии R 3.5.2, установленной в Microsoft Windows и Mac OS X, хотя они, вероятно, будут работать с любой текущей версией R.

Загрузите файлы примеров кода

Пакет с примерами кода для этой книги размещен в GitHub по адресу <https://github.com/PacktPublishing/Machine-Learning-with-R-Third-Edition> и по адресу <https://github.com/dataspelunking/MLwR/>.

Для того чтобы скачать файлы кода, нужно выполнить следующие действия.

1. Перейдите по указанной ссылке на сайт **github.com**.
2. Нажмите кнопку **Clone or Download**.
3. Щелкните кнопкой мыши на ссылке **Download ZIP**.
4. Скачайте архив с файлами примеров.

После загрузки файла распакуйте папку, используя последнюю версию одного из следующих архиваторов:

- WinRAR/7-Zip для Windows;
- Zipeg/iZip/UnRarX для Mac;
- 7-Zip/PeaZip для Linux.

Цветные иллюстрации

Мы также предоставляем PDF-файл с цветными скриншотами и схемами, приведенными в книге. Вы можете скачать его по адресу https://www.packtpub.com/sites/default/files/downloads/9781788295864_ColorImages.pdf.

Условные обозначения

В издании вы увидите несколько стилей текста, с помощью которых выделяются разные виды информации. Вот несколько примеров этих стилей и объяснение их значения.

Код в тексте, имена функций, имена файлов, расширения файлов, пользовательский ввод и названия R-пакетов отображаются следующим образом: «Функция `knn()` в пакете `class` предоставляет стандартную классическую реализацию алгоритма k-NN».

Пользовательский ввод и вывод в среде R записывается следующим образом:

```
> table(mushrooms$type)  edible
poisonous      4208      3916
```

Новые термины выделены курсивом, а **важные слова** — жирным шрифтом. Слова, которые вы видите на экране, например в меню или диалоговых окнах, выделены в тексте следующим образом: «Ссылка **Task Views** в левой части страницы **CRAN** указывает на список рекомендованных пакетов».



Важные примечания выглядят так.



Советы и подсказки описаны в таких врезках.

От издательства

Ваши замечания, предложения, вопросы отправляйте по адресу comp@piter.com (издательство «Питер», компьютерная редакция).

Мы будем рады узнать ваше мнение!

На веб-сайте издательства www.piter.com вы найдете подробную информацию о наших книгах.

1. Введение в машинное обучение

Если верить фантастике, то изобретение искусственного интеллекта неизбежно ведет к апокалиптическим войнам между машинами и их создателями. Эти истории начинаются с сегодняшней реальности: сначала мы обучаем компьютеры играть в простые игры, такие как крестики-нолики, и автоматизировать рутинные задачи. Затем, как пишут в этих книжках, машинам передают контроль над светофорами и коммуникациями, потом следуют военные беспилотники и ракеты. Эволюция машин принимает зловещий оборот, когда компьютеры становятся разумными и начинают обучать сами себя. Больше не нуждаясь в людях-программистах, они «удаляют» человечество.

К счастью, на момент написания этой книги машины все еще требовали ввода данных от пользователя.

Возможно, ваши впечатления от машинного обучения все еще до некоторой степени определяются картинками из средств массовой информации. Однако современные алгоритмы не являются универсальными и позволяют решать лишь конкретные задачи, чтобы представлять какую-либо опасность превратиться в самосознание. Цель современного машинного обучения — не создать искусственный мозг, а помочь нам разобраться в огромных запасах генерируемых человечеством данных.

К концу этой главы вы оставите в стороне популярные заблуждения и начнете лучше понимать суть машинного обучения. Вы также познакомитесь с основными понятиями, которые определяют и позволяют различать популярные подходы машинного обучения. Мы рассмотрим:

- происхождение, приложения и подводные камни машинного обучения;
- способы, которыми компьютеры превращают данные в знания и действия;
- этапы сопоставления алгоритмов машинного обучения с данными.

Технология машинного обучения предоставляет набор алгоритмов, преобразующих данные в руководство к действию. Продолжив читать книгу, вы увидите, как легко использовать R, чтобы начать применять машинное обучение для решения реальных задач.

Происхождение машинного обучения

Начиная с рождения, мы завалены данными. На датчики нашего тела — глаза, уши, нос, язык — постоянно воздействуют потоки необработанных данных, которые наш мозг преобразует в зрительные образы, звуки, запахи, вкусы. Используя речь, мы можем поделиться этим опытом с другими.

С появлением письменности человек начал записывать свои наблюдения. Охотники следили за движением стад животных; первые астрономы фиксировали расположение планет и звезд; в городах регистрировались налоговые поступления и т.д. Сегодня такие и многие другие наблюдения становятся все более автоматизированными и систематически регистрируются в постоянно растущих компьютеризированных базах данных.

Изобретение электронных датчиков способствовало новому взрывному росту объемов и разнообразия регистрируемых данных. Специальные датчики, такие как камеры, микрофоны, химические «носы», электронные «языки» и датчики давления, имитируют способность человека видеть, слышать, обонять, ощущать вкус и осязать. Эти датчики обрабатывают данные совсем не так, как это делает человек. В отличие от человека, отличающегося ограниченным и субъективным восприятием, электронный

датчик никогда не делает перерывы и не выражает эмоций, способных исказить его «ощущения».



Датчики видят объективную информацию, однако далеко не всегда предоставляют однозначное изображение реальности. Одни датчики имеют погрешность измерения вследствие аппаратных ограничений. У других ограничена область применения. Черно-белая фотография дает одно описание предмета, цветной снимок — другое. Микроскоп и телескоп позволяют увеличить изображение, но используются для объектов, расположенных на разных расстояниях.

В базах данных фиксируются многие показатели нашей жизни. Правительства, предприятия и частные лица регистрируют и используют всевозможную информацию — от глобально важной до бытовой. Датчики погоды записывают данные о температуре и давлении; камеры наблюдения следят за тротуарами и тоннелями метро; отслеживаются всевозможные детали цифрового следа: транзакции, переписка, связи в социальных сетях и многие другие.

Этот поток данных породил заявления о том, что мы вступили в эру *больших данных*, но это не совсем так. Люди всегда были окружены большими объемами данных. Уникальность нынешней эпохи в том, что теперь у нас есть огромное количество **зарегистрированных данных**, значительная часть которых напрямую доступна с компьютеров. Все больше интересных данных становятся доступными почти моментально, стоит только поискать в Интернете. Это огромное количество информации может побудить к обоснованным действиям, если систематически извлекать из нее осмысленные данные.

Область исследования, занимающаяся созданием компьютерных алгоритмов для преобразования данных в обоснованные действия, называется *машинным обучением*. Эта область возникла в среде, где доступные данные, статистические методы и вычислительные мощности развивались быстро и одновременно. Рост объемов данных потребовал дополнительных вычислительных мощностей, что, в свою очередь, стимулировало разработку статистических методов анализа больших наборов данных. Это породило цикл развития (схематично представлен на рис. 1.1), позволивший собирать все больше интересных данных и сформировавший современную среду, в которой практически по любой теме доступны бесконечные потоки данных.

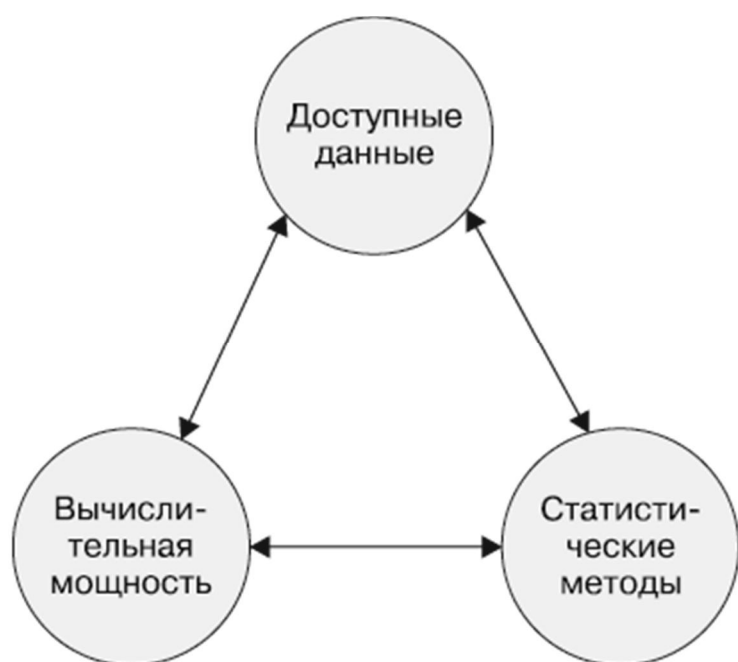


Рис. 1.1. Цикл развития, который сделал возможным машинное обучение

С машинным обучением тесно связан его родной брат — *интеллектуальный анализ данных* (data mining), который занимается генерированием новых идей из баз данных большого объема. Как следует из названия, интеллектуальный анализ данных включает в себя систематическую охоту на зачатки настоящего интеллекта. Несмотря на некоторые разногласия по поводу того, насколько сильно пересекаются машинное обучение и интеллектуальный анализ данных, потенциальное различие между ними заключается в том, что машинное обучение фокусируется на обучении компьютеров тому, как применять данные для решения задач, а интеллектуальный анализ данных направлен на обучение компьютеров выявлять закономерности, по которым люди обычно решают свои задачи.

Интеллектуальный анализ данных практически всегда подразумевает использование машинного обучения, но не любое машинное обучение требует интеллектуального анализа данных. Например, можно применить машинное обучение для сбора данных об автомобильном трафике для построения моделей, связанных с частотой аварий. С другой стороны, если компьютер учится водить автомобиль, то это чисто машинное обучение без интеллектуального анализа данных.



Выражение «интеллектуальный анализ данных» также иногда используется в уничижительном смысле для описания практики обманного сбора данных в поддержку некоей теории.

Область применения машинного обучения и злоупотребление им

Большинство людей слышали о Deep Blue — шахматном компьютере, который в 1997 году впервые выиграл у чемпиона мира. Другой известный компьютер, Watson, победил двух оппонентов в телевикторине Jeopardy (у

нас известна как «Своя игра») в 2011 году. Основываясь на этих потрясающих достижениях, эксперты предположили, что компьютерный интеллект заменит людей в сфере информационных технологий точно так же, как сейчас машины заменили рабочих на полях и сборочных линиях.

Правда заключается в том, что, даже когда машины достигают таких впечатляющих результатов, их способности полностью понять проблему заметно ограничены. Компьютер может лучше, чем человек, находить неявные закономерности в больших данных, но все равно требуется участие человека, чтобы проанализировать данные и превратить результат в осмысленное действие.



Мы не будем углубляться в обсуждение достижений Deep Blue и Watson, однако важно отметить, что ни один из этих компьютеров не умнее обычного пятилетнего ребенка. Подробнее о том, почему «меряться умами — дело скользкое», читайте в статье Уилла Грюнвальда (Will Grunewald) FYI: Which Computer Is Smarter, Watson Or Deep Blue?, опубликованной в журнале Popular Science в 2012 году (<https://www.popsci.com/science/article/2012-12/fyi-which-computer-smarter-watson-or-deep-blue>).

Машины не умеют задавать вопросы и даже не знают, какие вопросы следует задавать. Если вопрос сформулирован понятно для машины, то компьютер ответит на него гораздо лучше, чем человек. Современные алгоритмы машинного обучения сотрудничают с людьми так же, как гончие собаки — с охотниками: обоняние собаки во много раз сильнее, чем у ее хозяина, но без верного направления гончая может в итоге гоняться только за своим хвостом (рис. 1.2).

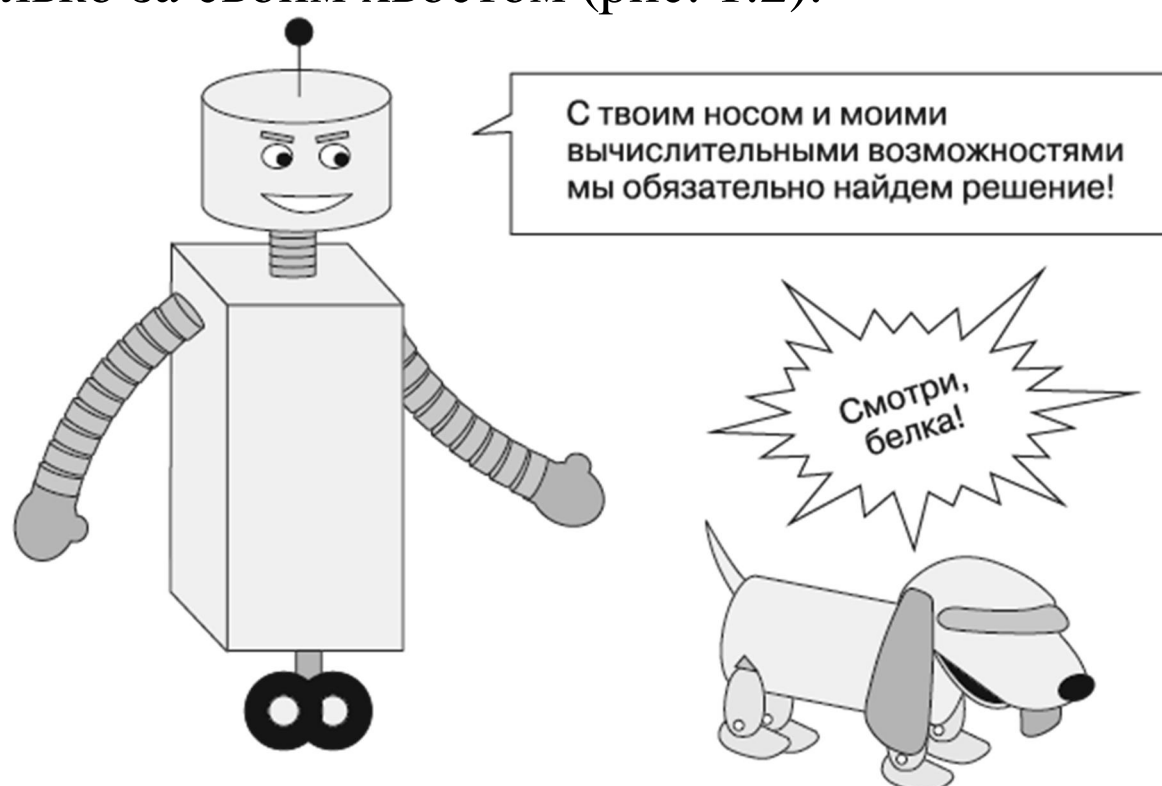


Рис. 1.2. Алгоритмы машинного обучения — мощные инструменты, которые требуют тщательного руководства

Чтобы лучше понять, где можно использовать машинное обучение, мы рассмотрим несколько примеров его успешного применения, случаи, где еще есть куда расти, и ряд ситуаций, в которых оно может принести больше вреда, чем пользы.

Успехи машинного обучения

Машинное обучение наиболее успешно там, где оно дополняет, а не заменяет специализированные знания эксперта в определенной области. Врачи применяют ML в борьбе с раком; оно помогает инженерам и программистам строить более умные дома и автомобили, а социологам — накапливать знания о том, как функционирует общество. С этой целью машинное обучение используется на многих предприятиях, в научных лабораториях, больницах и правительственных организациях. Везде, где производятся или накапливаются данные, скорее всего, используется хотя бы один алгоритм ML.

Невозможно перечислить все варианты использования машинного обучения, однако, рассмотрев недавние истории успеха, можно выделить несколько ярких примеров:

- выявление спама в электронной почте;
- разделение клиентов по поведению для создания таргетированной рекламы;
- прогнозы изменений погоды и долгосрочных изменений климата;
- борьба с мошенническими операциями по кредитным картам;
- страховые оценки финансового ущерба от штормов и стихийных бедствий;
- прогнозирование результатов всеобщих голосований;
- разработка алгоритмов автопилота для беспилотных летательных аппаратов и автомобилей;
- оптимизация расходования энергии в жилых и офисных зданиях;
- прогнозирование районов, где наиболее вероятна преступная деятельность;
- обнаружение цепочек генов, связанных с болезнями.

К концу книги вы познакомитесь с основными алгоритмами ML, которые используются для обучения компьютеров выполнению этих задач. Пока достаточно сказать, что, независимо от контекста, процесс машинного обучения одинаков. При любой задаче алгоритм берет данные и выявляет закономерности, которые позволяют принимать обоснованные решения.

Пределы возможностей машинного обучения

Несмотря на широкое применение и огромный потенциал машинного обучения, важно понимать его границы. В настоящее время машинное обучение имитирует относительно ограниченный набор возможностей человеческого мозга. Его гибкость для экстраполяции за пределами строгих параметров невелика, и ему не знаком здравый смысл. Имея это в

виду, следует проявлять предельную осторожность и точно определять, чему научился алгоритм, прежде чем применять его в реальном мире.

Не имея жизненного опыта, компьютеры также ограничены в способности делать простые выводы о следующих логических шагах. Возьмем, к примеру, рекламные баннеры, встречающиеся на многих сайтах. Они отображаются в соответствии с закономерностями, определенными с помощью интеллектуального анализа истории миллионов пользователей. Основываясь на этих данных, алгоритм делает вывод, что тот, кто просматривает сайты, продающие обувь, заинтересован в ее покупке и поэтому должен увидеть рекламу обуви. Проблема заключается в том, что это превращается в бесконечный цикл, в котором даже после покупки обуви пользователю предлагается дополнительная реклама обуви, а не шнурков для ботинок и крема для ухода за обувью.

Многие сталкивались с неполноценностью машинного обучения в понимании или переводе языка, а также распознаванию речи и рукописного текста. Возможно, самый первый пример подобного рода неудач — эпизод из телесериала «Симпсоны» (1994), в котором демонстрировалась пародия на планшет Apple Newton (рис. 1.3). В свое время Newton был известен передовой на то время системой распознавания рукописного текста. К сожалению для Apple, иногда эта система оказывалась неэффективной. В эпизоде телесериала это было обыграно так: Newton расшифровал хулиганскую записку «Избить Мартина» как «Съесть Марту».

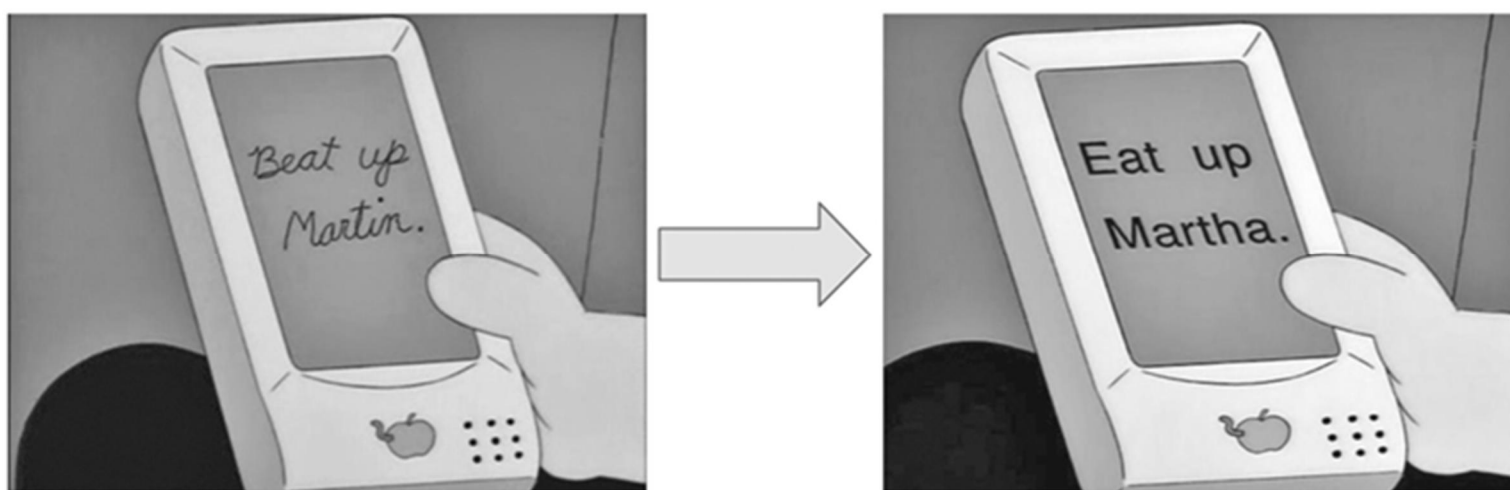


Рис. 1.3. Кадры из эпизода «Лиза на льду» из сериала «Симпсоны», 20th Century Fox (1994)

С тех пор машинная обработка естественного языка в Apple Newton значительно улучшилась, и теперь Google, Apple и Microsoft уверенно предлагают голосовые сервисы виртуального помощника, такие как Google Assistant, Siri и Cortana. Но и эти сервисы часто затрудняются ответить на относительно простые вопросы.

Кроме того, сервисы онлайн-перевода иногда неверно истолковывают предложения, которые даже ребенок легко понимает, а функция прогнозирования ввода текста на многих устройствах привела к появлению многочисленных юмористических сайтов с примерами «сбоя при

автозамене», которые иллюстрируют, что компьютеры способны понимать базовый язык, но совершенно не понимают контекст.

Некоторые из этих ошибок, безусловно, ожидаемы. Язык сложен, имеет смысловую многослойность и подтекст, и даже люди иногда неправильно понимают контекст. Несмотря на то что машинное обучение быстро совершенствуется в области обработки естественного языка, постоянно проявляющиеся недостатки демонстрируют тот важный факт, что оно хорошо ровно настолько, насколько хороши обработанные им данные. Если контекст во входных данных явно не представлен, то компьютеру, как и человеку, приходится делать лучшую догадку на основе своего ограниченного прошлого опыта.

Этика машинного обучения

По своей сути машинное обучение — просто инструмент, который помогает нам разобраться в сложных мировых данных. Как и любой инструмент, его можно использовать во благо или во зло. Хуже всего, когда машинное обучение применяется настолько широко или так жестко, что к людям относятся как к лабораторным крысам, автоматам или бездушным потребителям. На первый взгляд безвредный процесс, автоматизированный при помощи безэмоционального компьютера, может привести к непредвиденным последствиям. Именно поэтому было бы упущением, обсуждая использование машинного обучения и интеллектуальный анализ данных, хотя бы кратко не рассмотреть возможные этические последствия.

Из-за относительной молодости ML как дисциплины и той скорости, с которой оно развивается, сопутствующие правовые вопросы и социальные нормы часто бывают весьма неопределенны и постоянно меняются. Получая или анализируя данные, следует соблюдать осторожность, чтобы не нарушить закон, соблюсти условия обслуживания и соглашения об использовании данных, не злоупотребить доверием, не нарушить конфиденциальность клиентов или общественности.



Неофициальный корпоративный девиз Google — организации, которая собирает, возможно, больше данных о людях, чем кто-либо другой в мире, некогда звучал так: «Не причини зла» (англ. don't be evil). Все кажется вполне понятным, однако этого может оказаться недостаточно. Лучшим подходом может стать следование клятве Гиппократова, медицинскому принципу, который гласит: «Прежде всего — не навреди».

Компании розничной торговли обычно используют ML для создания таргетированной рекламы, складского учета и размещения товаров на полках магазинов. Многие из них снабжают свои кассы устройствами,

которые печатают купоны рекламных акций на основе истории покупок клиента. В обмен на небольшое количество личных данных клиент получает скидки на определенные продукты, которые он, вероятно, захочет купить. На первый взгляд это кажется относительно безвредным, однако посмотрим, что произойдет, если такая практика будет развиваться дальше.

Одна из историй, возможно выдуманных, касается крупной американской сети розничной торговли, в которой с помощью машинного обучения выявляли будущих мам, чтобы рассылать им купоны. Торговая сеть рассчитывала, что если будущим матерям предоставить значительные скидки, то они станут лояльными покупателями, которые впоследствии будут приобретать выгодные товары, такие как подгузники, детские смеси и игрушки.

Вооружившись методами машинного обучения, торговая сеть выделила в истории покупок товары, по которым можно с высокой степенью вероятности спрогнозировать не только то, является ли женщина беременной, но и приблизительные сроки родов.

После того как торговая сеть использовала эти данные для рекламных рассылок, некий человек связался с представителями сети и потребовал объяснить, почему его дочь получила купоны на товары для беременных. Он был в ярости от того, что продавец, похоже, поощрял беременность в подростковом возрасте! Как гласит история, когда из торговой сети позвонили, чтобы принести извинения, в итоге именно отец извинился после того, как поговорил с дочерью и узнал, что она действительно беременна!

Независимо от того, правда это или нет, урок, извлеченный из этой истории, заключается в том, что необходимо подключать здравый смысл, прежде чем слепо применять результаты анализа, полученного методами машинного обучения. Это особенно верно в тех случаях, когда речь идет о конфиденциальной информации, такой как данные о состоянии здоровья. Будь эта торговая сеть немного осторожнее, можно было предвидеть такой сценарий и быть осмотрительнее, выбирая способ применения данных, полученных по результатам анализа в процессе машинного обучения.



Подробнее о том, как розничные торговые сети используют машинное обучение для выявления беременностей, читайте в статье Чарльза Духигга (Charles Duhigg) *How Companies Learn Your Secrets* в журнале *New York Times Magazine* в 2012

году: <https://www.nytimes.com/2012/02/19/magazine/shopping-habits.html>.

Поскольку алгоритмы машинного обучения применяются все активнее, мы постоянно обнаруживаем, что компьютеры могут изучать некоторые неудачные варианты поведения человека. К сожалению, это может

проявляться в расовой и гендерной дискриминации и в появлении негативных стереотипов. Например, исследователи обнаружили, что сервис онлайн-рекламы Google показывает объявления о высокооплачиваемой работе чаще мужчинам, чем женщинам, а объявления о проверке криминального прошлого — чаще темнокожим, чем белым людям.

Эти типы ошибок распространились не только в Кремниевой долине. Разработанный Microsoft чат-бот Twitter быстро отключили после того, как он начал пропагандировать нацизм и антифеминизм. Часто алгоритмы, которые на первый взгляд кажутся «нейтральными по содержанию», спустя время начинают отражать убеждения большинства или доминирующие идеологии. Алгоритм, созданный Beauty.AI для отражения объективности человеческой красоты, вызвал споры, когда в числе победителей оказались почти только белые люди. Представьте себе, какие последствия имел бы этот алгоритм, если бы он был применен для распознавания лиц, подозреваемых в преступной деятельности!



Подробнее о реальных последствиях машинного обучения и дискриминации читайте в статье Клер Чейн Миллер (Claire Cain Miller) *When Algorithms Discriminate*, вышедшей в журнале *New York Times* в 2015 году: <https://www.nytimes.com/2015/07/10/upshot/when-algorithms-discriminate.html>.

Для того чтобы избежать дискриминации при применении алгоритмов ML, в некоторых юридических областях введены законы, запрещающие использование расовых, этнических, религиозных или других личных данных в коммерческих целях. Однако простого исключения таких данных из проекта может быть недостаточно, поскольку алгоритмы машинного обучения все равно могут непреднамеренно их различать. Если какая-то группа людей стремится жить в конкретном регионе, покупает определенные продукты или ведет себя таким образом, что можно однозначно идентифицировать их как группу, алгоритмы машинного обучения могут сделать вывод о личной информации на основании других факторов. В таких случаях может потребоваться *полностью* деидентифицировать этих людей, исключив любые *потенциально* идентифицирующие их данные в дополнение к уже закрытым статусам.

Помимо юридических последствий, ненадлежащее использование данных может помешать в достижении цели. Клиенты могут почувствовать себя некомфортно или испугаться, если какие-то стороны их жизни, которые они считают личными, станут достоянием общественности. В последние годы несколько известных веб-приложений пережили массовый

отток пользователей, которым показалось, что их эксплуатируют. Это произошло после того, как в приложениях изменились условия соглашений об обслуживании или выяснилось, что данные применялись в целях, на которые пользователи первоначально не рассчитывали. Поскольку состав конфиденциальной информации зависит от контекста, возрастной категории и места жительства, сложно установить границы надлежащего использования личных данных. Прежде чем начать проект, целесообразно рассмотреть возможные последствия, характерные для той или иной культуры, а также изучить актуальные законы, например недавно внедренные в Европейском союзе **Общие правила защиты данных** (General Data Protection Regulation, GDPR) и неизбежные изменения в политике конфиденциальности, которые за ними последуют.



Тот факт, что данные можно использовать с определенной целью, еще не означает, что это следует делать.

Наконец, важно отметить, что, по мере того как алгоритмы машинного обучения становятся все более важными в нашей повседневной жизни, у недобросовестных людей появляется соблазн воспользоваться ими в корыстных целях. Иногда злоумышленники хотят лишь сломать алгоритмы ради забавы или дурной славы, организовав что-то вроде «бомбардировки Google» — обмана системы ранжирования страниц Google методом толоки (англ. crowdsourced method). В других случаях результат может быть более драматичным. Актуальный пример — недавняя волна так называемых ложных новостей и вмешательство в выборы путем манипулирования алгоритмами, которые показывали рекламу и давали рекомендации в соответствии с личностными характеристиками пользователя. Чтобы избежать предоставления таких возможностей посторонним лицам, при создании систем машинного обучения важно учитывать, как на эти алгоритмы может повлиять определенный человек или группа людей [1].



Исследователь соцсетей danah boyd (именно так, строчными буквами) выступила с речью на конференции Strata Data Conference 2017 в Нью-Йорке, где обсуждалась важность усиления защиты алгоритмов машинного обучения от злоумышленников. Краткий вариант ее выступления вы найдете по адресу <https://points.datasociety.net/your-data-is-being-manipulated-a7e31a83577b>.

Последствия атак злоумышленников на алгоритмы машинного обучения могут быть и смертельными. Исследователи показали, что, создавая

«враждебную атаку», которая незаметно искажает дорожный знак с помощью тщательно подобранных граффити, злоумышленник может заставить беспилотный автомобиль неверно истолковать знак остановки, что, возможно, приведет к аварии со смертельным исходом. Несмотря на отсутствие злого умысла, ошибки, как в программном обеспечении, так и человеческие, уже привели к ряду несчастных случаев со смертельным исходом при внедрении технологии беспилотного управления транспортными средствами Uber и Tesla. С учетом этих примеров крайне важно, в том числе с этической точки зрения, чтобы специалисты по машинному обучению учитывали, как их алгоритмы могут быть использованы в реальном мире.

Как учатся машины

Формальное определение машинного обучения, приписываемое ученому, специалисту по информатике Тому М. Митчеллу (Tom M. Mitchell), гласит, что машина учится всякий раз, когда может использовать свой опыт, так что каждый следующий раз, когда ей приходится решать ту же задачу, она делает это эффективнее. Это определение интуитивно понятно, однако в нем полностью игнорируется процесс того, как именно опыт может быть преобразован в действия. И конечно, всегда легче сказать, чем действительно сделать!

Человеческий мозг способен учиться с рождения, а вот компьютеру следует четко указать условия. Несмотря на то что вовсе не обязательно знать теоретические основы обучения, все же это поможет нам понимать, различать и использовать на практике алгоритмы машинного обучения.



Сравнивая машинное обучение с обучением человека, вы, возможно, взглянете на собственный разум в ином свете.

Независимо от того, является ли обучаемый человеком или машиной, процесс обучения одинаков. Его можно разделить на четыре взаимосвязанные составляющие:

- *хранение данных* использует наблюдение, память и вспоминание, чтобы обеспечить основу в виде фактов для дальнейших рассуждений;
- *абстрагирование* подразумевает преобразование хранимых данных в более широкие представления и концепции;
- *обобщение* использует абстрактные данные для создания знаний и умозаключений, которые будут определять действия в новых контекстах;
- *оценка* обеспечивает механизм обратной связи для измерения полезности полученных знаний и информирования о потенциальных улучшениях.

На рис. 1.4 процесс обучения представлен в виде четырех отдельных составляющих, однако они организованы таким образом лишь в иллюстративных целях. На самом деле весь учебный процесс неразрывно связан. У людей обучение происходит на подсознательном уровне. Мы вспоминаем, обобщаем, делаем выводы и интуитивные прогнозы в пределах границ нашего разума, и поскольку этот процесс скрыт, то любые различия между людьми относятся к неопределенному понятию субъективности. Компьютеры, наоборот, выполняют эти процессы явно [2], и, поскольку весь процесс прозрачен, полученные знания можно исследовать, передавать, использовать в последующих действиях и рассматривать как науку о данных.

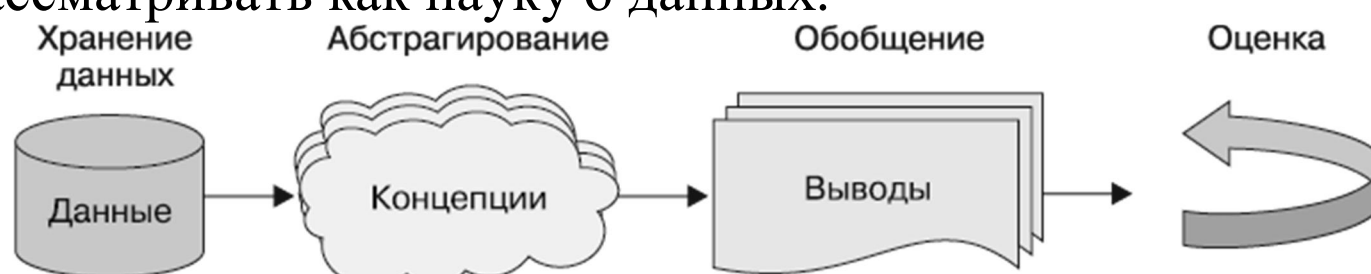


Рис. 1.4. Процесс обучения

Модное выражение «наука о данных» предполагает связь между данными, машиной и людьми, которые руководят процессом обучения. Все более широкое употребление этого термина в описаниях вакансий и различных академических трудах отражает его практическое применение как области исследований, связанной с теорией статистики и вычислений, а также с технологической инфраструктурой, обеспечивающей машинное обучение и сферы его применения. От специалистов в этой области требуется быть убедительными рассказчиками, сочетающими смелость в использовании данных с ограничениями того, на что могут повлиять эти данные и что можно прогнозировать с их помощью. Чтобы стать профессионалом в обработке данных, необходимо хорошо понимать, как работают алгоритмы обучения.

Хранение данных

Все обучение начинается с данных. Люди и компьютеры используют *хранение данных* как основу для более сложных рассуждений. Мозг человека принимает электрохимические сигналы, передаваемые в сети биологических клеток для хранения и обработки наблюдений в кратковременной и долговременной памяти. У компьютеров есть аналогичные возможности кратковременного и долговременного воспроизведения данных благодаря использованию жестких дисков, флеш-памяти и оперативной памяти (Random-Access Memory, RAM) в сочетании с центральным процессором (Central Processing Unit, CPU).

Тем не менее одной способности хранить и извлекать данные для обучения недостаточно. Сохраненные данные — это лишь единицы и нули на диске. Это набор воспоминаний, которые без более широкого контекста

являются бессмысленными. Без высокого уровня понимания знание — это всего лишь воспоминание, ограниченное тем, что было увидено ранее, и ничто иное.

Чтобы лучше понять суть этой идеи, вспомните, как вы в последний раз готовились к сложному тесту, например, перед выпускным экзаменом в университете, или к аттестации. Вы бы хотели иметь эйдетическую (фотографическую) память? Если это так, то вы, возможно, будете разочарованы, узнав, что идеальное воспроизведение вам вряд ли помогло бы. Даже если бы вы могли отлично запоминать материал, это бесполезно, если не знать точных вопросов и ответов, которые будут на экзамене. В противном случае пришлось бы запомнить ответы на каждый вопрос, который **можно было бы** задать по предмету, а вопросов — бесконечное количество. Очевидно, что это ненадежная стратегия.

Лучше потратить время на то, чтобы выборочно запомнить относительно небольшой набор репрезентативных идей, развивая при этом понимание того, как эти идеи взаимосвязаны и могут быть применены в случае непредвиденных обстоятельств. Таким образом определяются более широкие закономерности, и не требуется запоминать все детали, нюансы и потенциальные способы применения.

Абстрагирование

В процессе *абстрагирования* выполняется работа по выявлению более широкого значения хранимых данных, когда необработанные данные представляют более широкую, абстрактную концепцию или идею. Этот тип взаимосвязи, скажем, между объектом и его представлением, иллюстрируется знаменитой картиной Рене Магритта «Вероломство образов», представленной на рис. 1.5. На рисунке изображена курительная трубка с подписью *Ceci n'est pas une pipe* («Это не трубка»). Магритт продемонстрировал, что представление о трубке — это не сама трубка. Однако, несмотря на то что трубка ненастоящая, любой, кто смотрит на картину, видит именно трубку. Это говорит о том, что человек способен сопоставить **изображение** трубки с **идеей** трубки и с образом в памяти о **настоящей** трубке, которую можно держать в руке. Подобные абстрактные связи являются основой **представления знаний**, формирования логических структур, которые помогают превратить необработанную сенсорную информацию в осмысленное понимание.



Рис. 1.5. «Это не трубка». Источник: <http://collections.lacma.org/node/239578>

В процессе машинного представления знаний компьютер суммирует сохраненные необработанные данные, используя *модель* — явное описание закономерностей, связанных с данными. Подобно трубке Магритта, представление модели продолжает жить за пределами необработанных данных. Оно представляет идею большую, чем сумма ее частей.

Существует много разных типов моделей. Возможно, вы уже знакомы с некоторыми из них, например:

- математические уравнения;
- деревья и графы;
- логическое правило «если-то-иначе»;
- группы данных, известные как кластеры.

Выбор модели, как правило, не оставлен на усмотрение машины. Выбор определяют задача и тип имеющихся данных. О методах выбора подходящего типа модели речь пойдет далее в этой главе.

Процесс подгонки модели к набору данных называется *обучением*. После того как модель обучена, данные преобразуются в абстрактную форму, которая делает выводы из исходной информации.



Возможно, вы спросите, почему этот шаг называется обучением, а не учебой. Во-первых, обратите внимание, что процесс учебы не заканчивается абстрагированием данных — учащийся еще должен обобщить и оценить то, что было получено в процессе обучения. Во-вторых, слово «обучение» лучше отражает тот факт, что учитель-человек обучает ученика-машину понимать данные определенным образом.

Важно отметить, что обученная модель сама по себе не предоставляет новых данных, но в то же время ведет к получению новых знаний. Как это может быть? Ответ заключается в том, что наложение структуры на известные данные дает представление о том, чего мы еще не знаем. Это предполагает новое правило, по которому могут быть связаны элементы данных.

Рассмотрим, к примеру, открытие гравитации. Подбирая уравнения к данным наблюдений, Исаак Ньютон вывел концепцию гравитации, но сила, которую мы теперь называем гравитационной, существовала всегда. Ее просто не замечали, пока Ньютон не описал гравитацию как абстрактную концепцию, связывающую одни данные с другими, — а именно, введя константу g в модель, которая объясняет наблюдения за падающими объектами (рис. 1.6).



Рис. 1.6. Модели — это абстракции, которые объясняют наблюдаемые данные

Большинство моделей не приводит к развитию теорий, которые бы пошатнули научную мысль, сложившуюся на протяжении веков. Тем не менее в процессе абстрагирования можно выявить важные, но ранее не замечаемые закономерности и связи между данными. Модель, обученная на данных о геноме, может найти несколько генов, сочетание которых ответственно за возникновение диабета; банки могут обнаружить, казалось бы, безобидный тип транзакций, которые систематически появляются перед началом мошеннической деятельности; психологи могут определить комбинацию личностных характеристик, указывающих на какое-то расстройство. Эти закономерности были всегда, но, поскольку они представляют информацию в другом формате, из них можно сформулировать новую идею.

Обобщение

Следующим шагом в процессе обучения является использование абстрактных знаний для планирования действий. Однако среди бесчисленных базовых паттернов, которые могут быть идентифицированы в процессе абстрагирования, и бесконечных способов моделирования этих паттернов одни из них будут более полезными, чем другие. Если не ограничить создание абстракций неким полезным множеством, то обучаемый застрянет на том же месте, откуда начал, имея большой объем информации, но не зная, как дать ей практическую оценку.

Формально понятие «обобщение» означает процесс превращения абстрагированных знаний в форму, которую можно использовать для решения задач, подобных, но не идентичных тем, с которыми обучаемый уже сталкивался. Обобщение действует как поиск по всему набору

моделей (теорий или выводов), которые могут быть созданы на основе данных в процессе обучения.

Если бы мы могли представить себе множество, содержащее все возможные способы абстрагирования данных, то обобщение означало бы сокращение этого множества до меньшего, более управляемого множества важных выводов.

При обобщении перед обучаемым стоит задача ограничить обнаруживаемые им паттерны только теми, которые будут наиболее релевантными для решения последующих задач. Как правило, невозможно уменьшить количество паттернов, рассматривая их один за другим и ранжируя по степени полезности в будущем. Алгоритмы ML обычно используют ускоренные методы, которые быстрее сокращают пространство поиска. Для этого в алгоритмах применяются *эвристики*, которые представляют собой обоснованные предположения о том, где могут находиться наиболее полезные выводы.



В эвристиках используются аппроксимация и другие эмпирические правила. Это означает, что эвристики не гарантируют, что будет найдена наилучшая модель данных. Однако без таких ускоренных методов найти полезную информацию в большом наборе данных было бы невозможно.

Люди часто используют эвристики для быстрого обобщения опыта, чтобы задействовать его в новых обстоятельствах. Если вы когда-либо при необходимости быстро принять решение прислушивались к своему внутреннему голосу, прежде чем полностью оценить ситуацию, то вы интуитивно использовали ментальные эвристики.

Невероятная способность человека быстро принимать решения часто зависит не от логики, подобной компьютерной, а от эвристик, управляемых эмоциями. Иногда это может привести к нелогичным выводам. Например, людей, которые боятся летать самолетами, больше, чем тех, кто боится ездить на автомобилях, несмотря на то что, по статистике, автомобили более опасны. Это можно объяснить эвристикой доступности — тенденцией людей оценивать вероятность события по тому, насколько легко они могут вспомнить примеры. Аварии, связанные с авиаперелетами, широко освещаются. Поскольку это травмирующие события, о них, вероятно, будут помнить очень долго, тогда как автомобильные аварии едва заслуживают упоминания в новостях.

Безрассудство неверного применения эвристик характерно не только для людей. Эвристики, используемые алгоритмами машинного обучения, также иногда приводят к ошибочным выводам. Говорят, что алгоритм характеризуется *смещением* (*bias*, можно перевести и как «предвзятость»), если его выводы *систематически* оказываются ошибочными. Это

предполагает, что выводы являются ошибочными последовательно или предсказуемо.

Предположим, что алгоритм машинного обучения научился распознавать лица, выделяя два темных круга, обозначающих глаза, которые расположены над прямой линией — ртом. В этом случае алгоритму сложно будет распознать лицо или же он может систематически ошибаться в отношении лиц, которые не соответствуют его модели. Алгоритм может не распознавать лица в очках, а также лица, повернутые под углом, смотрящие вбок или лица с определенными оттенками кожи, формой или особенностями, которые не соответствуют его восприятию. Схематично этот процесс изображен на рис. 1.7.

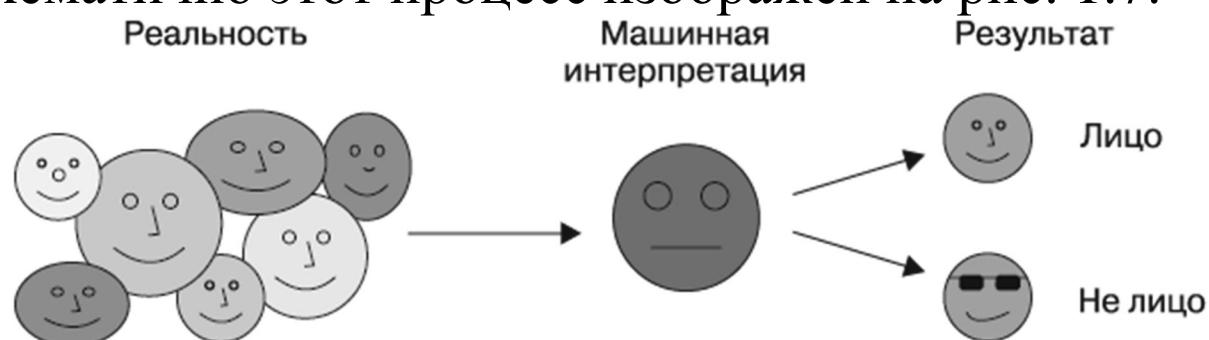


Рис. 1.7. В процессе обобщения опыта появляется смещение (систематическая ошибка)

В современном понимании «предвзятость» имеет весьма негативные коннотации. СМИ постоянно заявляют, что они свободны от предвзятости, и утверждают, что подают факты объективно, без эмоций. Тем не менее представьте себе, что незначительная предвзятость может быть полезной. Без некой доли своеволия бывает трудно выбрать один из нескольких конкурирующих вариантов, у каждого из которых есть свои плюсы и минусы, не так ли? Действительно, исследования в области психологии показали, что люди, рожденные с повреждениями участков головного мозга, ответственных за эмоции, могут испытывать сложности при принятии решений. Они способны часами обсуждать простые решения — например, какого цвета рубашку надеть или куда пойти на обед. Как это ни парадоксально, именно предвзятость ограждает нас от части лишней информации, а также позволяет использовать другую информацию для действий. Именно так алгоритмы машинного обучения выбирают способ понимания данных из бесчисленного множества, предлагаемого им для изучения.

Оценка

Предвзятость — неизбежное зло, связанное с процессами абстракции и обобщения, характерными для любого обучения. Для того чтобы выбрать действие в условиях безграничных возможностей, необходима предвзятость на протяжении всего обучения. Следовательно, у каждой стратегии обучения есть свои недостатки; не существует единого алгоритма обучения. Таким образом, последний этап учебного процесса —

оценка его успешности и определение степени успеваемости ученика, несмотря на предвзятость. Информацию, полученную на этапе оценки, затем можно использовать для принятия решения о необходимости дополнительного обучения.



После того как вы успешно освоите какую-нибудь из техник машинного обучения, у вас может возникнуть желание применять ее ко всем задачам. Главное — не поддаваться этому искушению, поскольку ни один подход ML не является универсальным для всех возможных ситуаций. Это обстоятельство описывается «теоремой об отсутствии бесплатных завтраков» (No Free Lunch), сформулированной Дэвидом Вольпертом (David Wolpert) в 1996 году. Подробнее об этом читайте на сайте <http://www.no-free-lunch.org>.

Как правило, оценка выполняется после окончания обучения модели на начальном *тренировочном наборе данных (training dataset)*. Затем модель оценивается по отдельному *тестовому набору данных (test dataset)*, чтобы определить, насколько хорошо характеристика тренировочных данных позволяет делать обобщения для новых данных, не входящих в тренировочный набор (unseen cases). Следует отметить, что модель крайне редко позволяет делать обобщения для данных, не входящих в тренировочный набор, — ошибки почти всегда неизбежны.

В частности, модели не позволяют делать идеальные обобщения из-за *шума (noise)* — термин описывает необъясненные или необъяснимые изменения данных. Зашумление данных может быть вызвано внешне случайными событиями, такими как:

- ошибка измерения из-за неточности датчиков, которые иногда прибавляют к показаниям небольшие значения (или вычитают их);
- проблемы, связанные с поведением человека, например, когда респонденты дают случайные ответы на вопросы, чтобы быстрее закончить опрос;
- проблемы качества данных, в том числе отсутствующие, неполные, неправильно закодированные или искаженные значения;
- процессы настолько сложные или малопонятные, что влияют на данные способами, которые кажутся случайными.

Попытка смоделировать шум лежит в основе проблемы, называемой *переобучением*, или *перетренировкой (overfitting)*. Поскольку большинство зашумленных данных необъяснимы по определению, попытки объяснить шум приводят к появлению моделей, которые не будут хорошо обобщаться для новых случаев, а также более сложных моделей, не соответствующих реальному паттерну, который обучаемый пытается выявить (рис. 1.8).

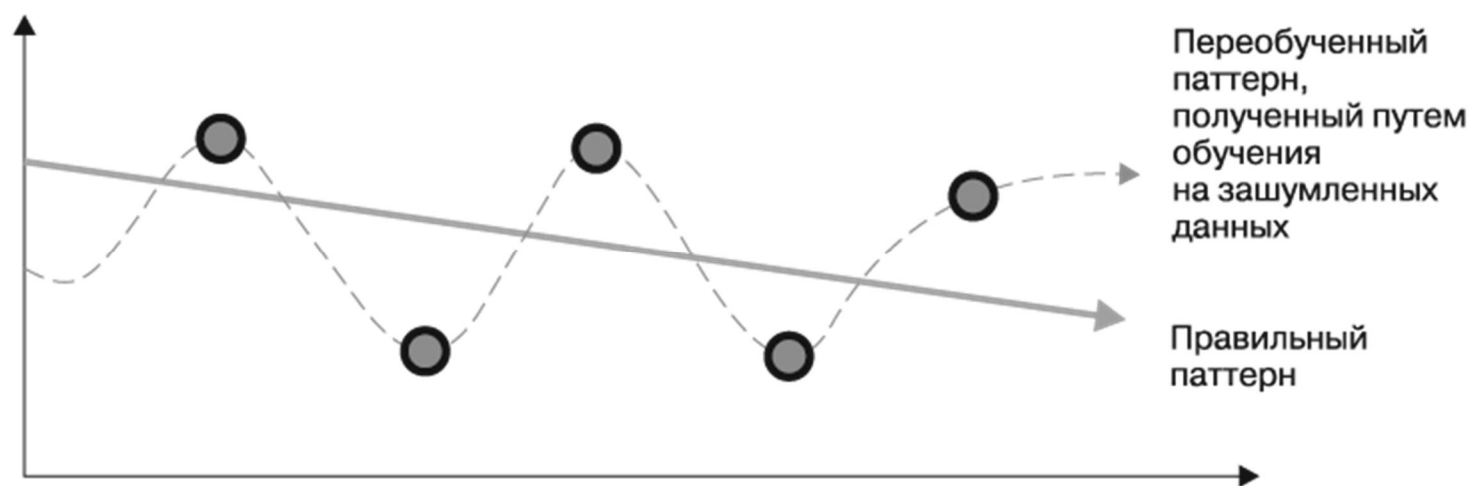


Рис. 1.8. Моделирование шума обычно приводит к появлению более сложных моделей, в которых упускаются основные паттерны

Считается, что модель, которая сравнительно хорошо работает на этапе обучения, но сравнительно плохо — на этапе оценки, является *переобученной* для тренировочного набора данных, потому что плохо поддается обобщению для тестового набора данных. С практической точки зрения это означает, что модель выявила закономерность в данных, бесполезную для предстоящих действий; процесс обобщения не удался. Для каждого конкретного подхода машинного обучения применяются свои решения проблемы перетренированности. Пока нам нужно лишь знать об этой проблеме. Важным показателем качества методов обучения является то, насколько хорошо они способны обрабатывать зашумленные данные и избегать перетренированности.

Машинное обучение на практике

До сих пор речь шла о том, как машинное обучение работает в теории. Чтобы применить процесс обучения к реальным задачам, мы будем использовать пятиэтапный процесс. Независимо от задачи, любой алгоритм ML можно разбить на следующие этапы.

1. Сбор данных. Этап включает в себя сбор материалов для обучения, которые алгоритм будет использовать для получения практических знаний. В большинстве случаев данные необходимо собрать в единый источник, такой как текстовый файл, электронная таблица или база данных.

2. Исследование и подготовка данных. Качество любого проекта по машинному обучению во многом определяется качеством входных данных. Таким образом, на этом этапе важно как можно больше узнать о данных и их особенностях. Необходимо выполнить дополнительную работу по подготовке данных к процессу обучения. Сюда входит исправление или очистка так называемых грязных данных, удаление ненужных данных и их перекодирование в соответствии с тем, какие входные данные ожидает получить обучаемый.

3. Обучение модели. К тому времени, как данные будут готовы для анализа, у вас, вероятно, уже будет представление о том, что можно извлечь из них. Выбор алгоритма определяется конкретной задачей

машинного обучения, а сам алгоритм предоставляет данные в форме модели.

4. Оценка модели. Каждая модель машинного обучения приводит к предвзятому решению задачи обучения. Это означает, что важно оценить, насколько хорошо алгоритм извлек уроки из своего опыта. В зависимости от типа используемой модели можно оценить ее точность с помощью тестового набора данных или же может потребоваться подобрать показатели производительности, специфичные для предполагаемой области применения.

5. Улучшение модели. Если нужна более высокая производительность, необходимо использовать улучшенные стратегии. Иногда может потребоваться вообще перейти на другой тип модели. Возможно, вам придется дополнить данные или выполнить подготовительную работу, как на втором этапе этого процесса.

После того как эти шаги выполнены, модель, если она, как вам кажется, работает хорошо, можно применить для решения поставленной задачи. В зависимости от обстоятельств можно использовать ее для предоставления расчетных данных прогнозов (возможно, в режиме реального времени); прогнозирования финансовых данных; проведения аналитического обзора для маркетинга; автоматизации некоторых задач. Успехи и неудачи используемой модели могут также предоставить дополнительные данные для последующего обучения.

Типы входных данных

Практика машинного обучения включает в себя сопоставление характеристик входных данных с предвзятостью имеющихся обучающих алгоритмов. Таким образом, прежде, чем применить ML для решения реальных задач, важно понимать терминологию, которая различает разные входные наборы данных.

Выражение «*единица наблюдения*» используется для описания наименьшего объекта с измеренными свойствами, представляющего интерес для исследований. Обычно единицей наблюдения выступают люди, объекты, транзакции, моменты времени, географические регионы или измерения. Иногда единицы наблюдения объединяются в сложные единицы, например человеко-года (когда наблюдение ведется за одним и тем же человеком в течение нескольких лет; каждый человеко-год включает в себя данные о человеке за один год).



Единица наблюдения связана с *единицей анализа*, но не идентична ей. Единица анализа — наименьшая единица, на основании которой делается вывод. Единицы наблюдения и анализа часто совпадают, но не всегда.

Например, данные, полученные от людей (единица наблюдения), можно использовать для анализа тенденций в разных странах (единица анализа).

Наборы данных, в которых хранятся единицы наблюдения и их свойства, можно объединить в коллекции:

- *примеров* — экземпляров единиц наблюдения, для которых были зарегистрированы свойства;
- *признаков* — зарегистрированных свойств или атрибутов примеров, которые могут пригодиться для обучения.

Понять, что такое признаки и примеры, проще всего на реальных сценариях. Так, для создания обучающего алгоритма, распознающего спам, единицами наблюдения могут выступать сообщения электронной почты, примерами — конкретные сообщения, а признаками могут выступать слова, используемые в сообщениях.

Для алгоритма диагностики рака единицей наблюдения могут быть пациенты, примеры могут включать в себя случайную выборку онкологических пациентов, а признаками могут быть геномные маркеры из клеток, полученных после биопсии, а также характеристики пациента, такие как вес, рост или артериальное давление.

Люди и машины различаются по типам сложности входных данных, которые они способны обработать. Людям удобно принимать *неструктурированные данные*, такие как произвольный текст, изображения или звук. Люди также гибко относятся к таким случаям обработки, в которых у одних наблюдений множество признаков, а у других — совсем мало. Компьютеры, напротив, обычно требуют *структурированных* данных: каждый пример явления имеет одинаковые признаки, которые представлены в форме, понятной компьютеру. Для использования машинного прямого перебора для больших неструктурированных наборов данных обычно требуется преобразовать входные данные в структурированную форму.

На рис. 1.9 данные представлены в *матричном формате*, где каждая строка таблицы является примером, а каждый столбец — признаком. В данном случае строки — это примеры продаваемых автомобилей, а в столбцах отражены характеристики каждого автомобиля, такие как цена, пробег, цвет и тип трансмиссии. Данные в матричном формате сегодня являются наиболее распространенной формой представления, используемой в машинном обучении. Как будет видно в следующих главах, где рассматриваются форматы данных в специализированных приложениях, перед началом машинного обучения эти форматы в итоге преобразуются в матрицу.

Признаки					
Год выпуска	Модель	Цена	Пробег	Цвет	Тип трансмиссии
2011	SEL	21 992	7413	Yellow	AUTO
2011	SEL	20 995	10 926	Gray	AUTO
2011	SEL	19 995	7351	Silver	AUTO
2011	SEL	17 809	11 613	Gray	AUTO
2012	SE	17 500	8367	White	MANUAL
2010	SEL	17 495	25 125	Silver	AUTO
2011	SEL	17 000	27 393	Blue	AUTO
2010	SEL	16 995	21 026	Silver	AUTO
2011	SES	16 995	32 655	Silver	AUTO

Примеры

Рис. 1.9. Простой набор данных, представленных в матричном формате: описание автомобилей, выставленных на продажу

Признаки набора данных могут быть представлены в различной форме. Если признак обозначает характеристику, измеряемую в числах, логично, что он именуется *числовым*. И напротив, если признак — набор категорий, то он называется *категориальным* или *номинальным*. Особая разновидность категориальной переменной называется *порядковой*, что означает номинальную переменную, категории которой образуют упорядоченный список. В качестве примеров порядковых переменных можно привести размеры одежды: маленький, средний и большой; определение степени удовлетворенности клиентов по шкале от «совершенно недоволен» до «не очень доволен» и «очень доволен». Какова бы ни была задача обучения, рассуждения о том, что представляют собой признаки ее набора данных, их типы и единицы измерения, помогут выбрать подходящий алгоритм машинного обучения.

Типы алгоритмов машинного обучения

Алгоритмы машинного обучения делятся на категории в соответствии с их назначением. Понимание категорий алгоритмов обучения является первым шагом к использованию данных для совершения желаемых действий. *Модель прогнозирования* применяется для задач, которые включают в себя, как следует из названия, прогноз одного значения с использованием других значений из набора данных. Обучающий алгоритм пытается обнаружить и смоделировать взаимосвязь между *целевым признаком* (прогнозируемым) и другими признаками.

В отличие от обычного трактования слова «прогноз» как предсказания модели прогнозирования не обязательно должны предвидеть предстоящие события. Например, модель прогнозирования можно использовать для предсказания прошлых событий, таких как дата зачатия ребенка на основании текущих уровней гормонов у матери. Модели прогнозирования также применяют для управления светофорами в час пик в режиме реального времени.

Поскольку модели прогнозирования получают четкие инструкции о том, что и как именно им нужно изучить, процесс обучения прогнозирующей модели называется *обучением с учителем*. Под учителем здесь подразумевается не участие человека, а скорее тот факт, что целевые значения дают обучаемому возможность определить, насколько хорошо он усвоил задачу. С более формальной точки зрения алгоритм обучения с учителем пытается оптимизировать функцию (модель), чтобы найти сочетание значений признаков, которые приводят к целевому результату для конкретного набора данных.

Часто встречается задача машинного обучения с учителем, состоящая в прогнозировании того, к какой категории относится пример. Такая задача называется *классификацией*. Представить себе потенциальное использование классификатора легко. Например, можно спрогнозировать следующие события:

- сообщение, полученное по электронной почте, является спамом;
- у этого пациента рак;
- футбольная команда выиграет;
- заявитель не выполнит своих обязательств по кредиту.

В задаче классификации целевой признак, который должен быть спрогнозирован, является категориальным признаком и называется *классом*. Класс делится на категории — *уровни*. Класс может иметь два уровня и более, а уровни могут быть упорядоченными или неупорядоченными. Классификация используется в машинном обучении так широко, что существует множество типов алгоритмов классификации. У каждого из них есть плюсы и минусы, каждый из них подходит для определенных типов входных данных. Примеры алгоритмов мы увидим и в этой главе, и в книге далее.

Алгоритмы, обучаемые с учителем, могут также использоваться для прогнозирования числовых данных, таких как доход, лабораторные показатели, результаты тестов или количество предметов. Для прогнозирования таких числовых значений предназначен популярный тип алгоритмов *числового прогнозирования*, соответствующих моделям линейной регрессии для входных данных. Хотя регрессия не единственный метод числового прогнозирования, она, безусловно, наиболее популярна. Методы регрессии широко используются для прогнозирования, поскольку они точно, количественно определяют взаимосвязь между исходными данными и целевым значением, включая как величину, так и степень неопределенности этой взаимосвязи.



Поскольку числа легко преобразовать в категории (например, возраст от 13 до 19 лет — подростки), а категории — в числа (например, 1 — мужчины, 0 — женщины), то граница между классификационными моделями и моделями числового прогнозирования не всегда является жесткой.

Описательная модель используется для задач, где можно извлечь выгоду из результатов обобщения данных, выполненного новыми, более интересными способами. В отличие от прогнозирующих моделей, которые предсказывают интересующий объект, в описательной модели все признаки одинаково важны. В сущности, поскольку цель обучения отсутствует, процесс обучения описательной модели называется *обучением без учителя*. Возможно, представить применение описательных моделей сложнее — в конце концов, что хорошего в том, что обучаемый ничего не изучает, — однако эти модели регулярно используются для интеллектуального анализа данных.

Например, задача описательного моделирования, называемая *обнаружением закономерности*, используется для выделения среди данных полезных ассоциаций. Обнаружение закономерностей — цель *анализа потребительской корзины*, который применяется к данным о транзакциях, касающимся розничных продаж. В этом случае розничные продавцы стремятся выделить сопутствующие товары, так что полученная информация может быть использована для совершенствования тактики маркетинга. Например, если продавец заметит, что плавки обычно покупают вместе с солнцезащитным кремом, то он может разместить эти предметы в магазине рядом или устроить акцию, чтобы продать их покупателям вместе.



Первоначально использовавшееся только в контексте розничной торговли, «обнаружение закономерностей» теперь стали задействовать весьма инновационным образом. Например, его можно применять для выявления моделей мошеннического поведения, генетических дефектов или «горячих точек» преступной деятельности.

Задача описательного моделирования по разделению набора данных на однородные группы называется *кластеризацией*. Она иногда применяется для *сегментного анализа рынка*, который позволяет идентифицировать группы людей с похожим поведением или по демографическим признакам, чтобы, основываясь на общих характеристиках, нацелить на них рекламные кампании. При таком подходе машина идентифицирует кластеры, но для их интерпретации требуется вмешательство человека. Например, если клиенты продуктового магазина разделились на пять кластеров, то отдел маркетинга должен будет понять различия между

этими группами, чтобы составить рекламную кампанию, которая лучше всего будет подходить для каждой из групп. Несмотря на необходимые усилия человека, это все равно требует меньше трудозатрат, чем создание уникальной привлекательности магазина для каждого клиента.

Наконец, класс алгоритмов машинного обучения, называемых *метаобучением*, не привязан к конкретной задаче обучения, а скорее сосредоточен на обучении тому, как учиться более эффективно. Алгоритм метаобучения использует результаты прошлого обучения для дополнительного обучения.

В эту группу входят алгоритмы обучения, которые учатся работать совместно, в группах, именуемых *ансамблями*, а также алгоритмы, которые, похоже, со временем, превращаются в процесс, называемый *обучением с подкреплением*. Метаобучение может быть полезным для очень сложных задач или в тех случаях, когда результат алгоритма прогнозирования должен быть максимально точным.

Одно из самых захватывающих исследований, выполняемых сегодня в области машинного обучения, относится к области метаобучения. Например, *сопоставительное обучение* включает в себя изучение слабых сторон модели для повышения ее продуктивности или защиты от вредоносных атак. Значительные инвестиции также вкладываются в исследования и разработки, направленные на создание более крупных и быстрых ансамблей, способных моделировать огромные наборы данных с использованием высокопроизводительных компьютеров или облачных сред.

Подбор алгоритмов по входным данным

В табл. 1.1 перечислены основные типы алгоритмов машинного обучения, описанные в данной книге. Это лишь часть всего множества алгоритмов машинного обучения, однако изучения перечисленных методов будет достаточно для понимания любых других методов, с которыми вы можете столкнуться.

Таблица 1.1. Основные типы алгоритмов машинного обучения

Модель	Задача обучения	Где прочитать
Алгоритмы обучения с учителем		
Метод k-ближайших соседей	Классификация	Глава 3
Наивный байесовский классификатор	Классификация	Глава 4
Деревья решений	Классификация	Глава 5
Правила классификации	Классификация	Глава 5
Линейная регрессия	Числовое прогнозирование	Глава 6
Регрессионные деревья	Числовое прогнозирование	Глава 6
Деревья моделей	Числовое прогнозирование	Глава 6
Нейронные сети	Двойное назначение	Глава 7
Метод опорных векторов	Двойное назначение	Глава 7
Алгоритмы обучения без учителя		
Ассоциативные правила	Обнаружение закономерностей	Глава 8
Кластеризация методом k-средних	Кластеризация	Глава 9
Алгоритмы метаобучения		
Бэггинг	Двойное назначение	Глава 11
Бустинг	Двойное назначение	Глава 11
Случайный лес	Двойное назначение	Глава 11

Для того чтобы использовать машинное обучение в реальном проекте, вам необходимо определить, какая из четырех задач обучения лежит в основе вашего проекта: классификация, числовое прогнозирование, обнаружение закономерностей или кластеризация. От этой задачи и будет зависеть выбор алгоритма. Например, если требуется обнаружение закономерностей, то вы, скорее всего, примените ассоциативные правила. Если же задача требует кластеризации, то, вероятно, будет использоваться алгоритм k-средних, а для числового прогнозирования — регрессионный анализ или регрессионные деревья.

Для классификации нужно потратить дополнительное время на исследование, чтобы выбрать соответствующий классификатор для данной задачи обучения. В таких случаях полезно рассмотреть различия между алгоритмами, которые станут очевидными только при углубленном изучении каждого классификатора. Например, в рамках задач классификации деревья решений позволяют построить понятные модели, в то время как модели нейронных сетей, как известно, трудны для интерпретации. При разработке модели кредитоспособности это может быть важным моментом, поскольку закон часто требует, чтобы заявитель был уведомлен о причинах отказа в выдаче ему кредита. Даже если нейронная сеть лучше прогнозирует вероятность неуплаты по кредиту, но при этом ее прогнозы невозможно объяснить, в данном случае она бесполезна.

Чтобы вам легче было определиться в выборе алгоритма, в каждой главе перечислены основные плюсы и минусы конкретного алгоритма обучения. Вы увидите, что иногда эти характеристики позволяют однозначно исключить из рассмотрения некоторые модели, однако во многих случаях выбор алгоритма является произвольным. Если это произойдет, не бойтесь

использовать любой алгоритм, который считаете более удобным. В других случаях, когда главная цель — точность прогнозирования, может потребоваться протестировать несколько моделей и выбрать ту из них, которая подходит лучше всего, или же использовать алгоритм метаобучения, который объединяет в себе несколько разных обучаемых моделей, чтобы задействовать сильные стороны каждой из них.

Машинное обучение с использованием R

Многие алгоритмы, необходимые для машинного обучения, не входят в базовую инсталляцию R. Однако они доступны благодаря обширному сообществу экспертов, которые свободно делятся результатами своей работы. Эти алгоритмы следует установить вручную поверх базовой инсталляции R. Благодаря тому что R является бесплатным программным обеспечением с открытым исходным кодом, вам не придется платить за эту дополнительную функциональность.

Набор функций R, свободно распространяемых между пользователями, называется *пакетом*. Для каждого алгоритма ML, описанного в книге, существуют бесплатные пакеты. В сущности, данная книга охватывает лишь небольшую часть существующих пакетов машинного обучения на R.

Если вас интересует более подробная информация об R-пакетах, обратитесь к списку *Comprehensive R Archive Network (CRAN)* — это всемирная коллекция сайтов и FTP-ресурсов, где представлены последние версии программного обеспечения и пакеты для R. Если вы получили программное обеспечение R путем скачивания, то, скорее всего, оно скачано с сайта CRAN, который доступен по адресу <http://cran.r-project.org/index.html>.



Если у вас еще не установлен R, то на сайте CRAN вы также найдете инструкции по его установке и информацию о том, куда обратиться при возникновении проблем.

Ссылка **Packages** в левой части веб-страницы CRAN ведет на страницу, где можно просмотреть список пакетов в алфавитном порядке или отсортировать их по дате публикации. На момент написания этой книги было доступно в общей сложности 13 904 пакета — вдвое больше, чем на момент написания второго издания, и более чем втрое — со времени выхода первого издания! Очевидно, что R-сообщество процветает, и эта тенденция не замедляется!

Ссылка **Task Views**, расположенная в левой части веб-страницы CRAN, ведет к списку пакетов, отсортированному по темам. Описание задач машинного обучения с перечислением пакетов, представленных в этой

книге (и многих других), вы найдете по адресу <https://CRAN.R-project.org/view=MachineLearning>.

Установка R-пакетов

Несмотря на обширный набор доступных расширений R, благодаря формату пакета его установка и использование практически не требуют усилий. Чтобы продемонстрировать использование пакетов, мы установим и загрузим пакет `RWeka`, разработанный Куртом Хорником (Kurt Hornik), Кристианом Бухтой (Christian Buchta) и Акимом Зайлисом (Achim Zeileis), — подробнее об этом пакете читайте в публикации *Open-Source Machine Learning: R Meets Weka, Computational Statistics, Vol. 24*, p. 225–232.

Пакет `RWeka` предоставляет набор функций, которые открывают для R доступ к алгоритмам машинного обучения, входящим в состав программного пакета `Weka`, написанного на Java и разработанного Иеном Х. Уиттеном (Ian H. Witten) и Эйбом Франком (Eibe Frank). Подробнее о пакете `Weka` читайте здесь: <http://www.cs.waikato.ac.nz/~ml/weka/>.



Для того чтобы использовать пакет `RWeka`, вам нужно установить Java, если это еще не сделано (на многих компьютерах Java устанавливается по умолчанию). Java — это бесплатный набор инструментов программирования, который позволяет использовать кросс-платформенные приложения, такие как `Weka`. Подробнее о Java читайте здесь: <http://www.java.com>. Там же вы сможете загрузить Java для вашей системы.

Самый прямой способ установки пакета — с помощью функции `install.packages()`. Для того чтобы установить пакет `RWeka`, просто введите в командной строке R следующую команду:

```
> install.packages("RWeka")
```

Среда R подключится к CRAN и загрузит пакет в формате, соответствующем вашей операционной системе. Некоторые пакеты, такие как `RWeka`, прежде чем их можно будет использовать, требуют установки дополнительных пакетов. Это так называемые *зависимости*. По умолчанию установщик автоматически загрузит и установит все зависимости.



При первой установке R-пакета установщик может попросить вас выбрать зеркало CRAN. В таком случае выберите ближайшее от вас зеркало. Этим вы обеспечите самую высокую скорость загрузки.

Параметры установки, предлагаемые по умолчанию, подходят для большинства систем. Однако иногда требуется установить пакет в другое место. Например, если у вас нет полномочий пользователя **root** или администратора, то вам может потребоваться выбрать другой путь установки. Это можно сделать с помощью параметра `lib` следующим образом:

```
> install.packages("RWeka", lib =  
"/path/to/library")
```

Функция установки также предоставляет дополнительные параметры для установки из локального файла, из заданного источника или для использования экспериментальных версий. Чтобы узнать больше об этих параметрах в файле справки, можно ввести следующую команду:

```
> ?install.packages
```

Вообще, вопросительный знак может использоваться для получения справки о любой R-функции. Для этого просто введите `?`, а затем название функции.

Загрузка и выгрузка R-пакетов

Для того чтобы сэкономить память, R не загружает по умолчанию все установленные пакеты. Их загружают пользователи по мере необходимости с помощью функции `library()`.



Из-за названия функции `library()` некоторые ошибочно считают, что термины «библиотека» и «пакет» являются взаимозаменяемыми. Однако, строго говоря, библиотека — это место, где установлены пакеты, а не сам пакет.

Для того чтобы загрузить ранее установленный пакет `RWeka`, введите следующую команду:

```
> library(RWeka)
```

Кроме `RWeka`, в следующих главах будут использоваться еще несколько R-пакетов. Инструкции по установке этих дополнительных пакетов будут приводиться по мере необходимости.

Для того чтобы выгрузить R-пакет, используется функция `detach()`. Например, чтобы выгрузить пакет `RWeka`, введите следующую команду:

```
> detach("package:RWeka", unload = TRUE)
```

При этом освободятся все ресурсы, используемые пакетом.

Установка RStudio

Перед началом работы с R настоятельно рекомендую также установить приложение *RStudio* с открытым исходным кодом. RStudio — это

дополнительный интерфейс к R, включающий в себя функциональные возможности, которые значительно упрощают работу с R-кодом, делают ее более удобной и интерактивной (рис. 1.10). Приложение RStudio доступно бесплатно по адресу <https://www.rstudio.com/>.

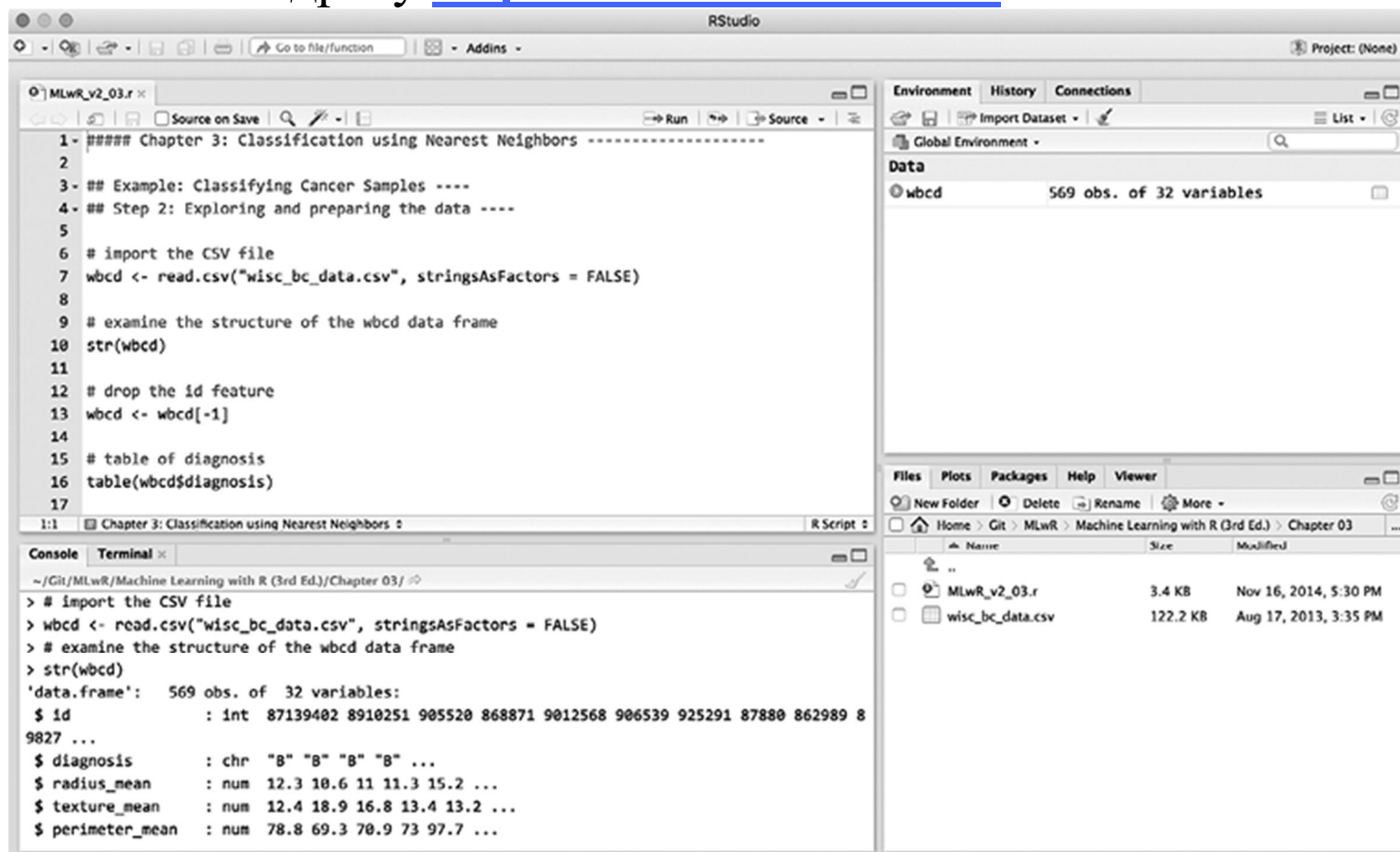


Рис. 1.10. Графическая среда RStudio делает использование R проще и удобнее

Интерфейс RStudio включает в себя встроенный редактор кода, консоль командной строки R, браузер файлов и браузер объектов R. Синтаксические конструкции кода R автоматически выделяются цветом, а результаты выполнения кода, диаграммы и графики отображаются непосредственно в среде, что значительно упрощает отслеживание длинных или сложных операторов и программ. Более продвинутые функции позволяют управлять R-проектами и пакетами; обеспечивают интеграцию с инструментами контроля исходного кода и контроля версий, такими как *Git* и *Subversion*; помогают управлять подключением к базе данных; обеспечивают преобразование результатов работы R-кода в форматы HTML, PDF и Microsoft Word.

Наличие RStudio — главная причина, по которой R сегодня является лучшим выбором для специалистов по обработке данных. Эта среда объединяет в себе возможности программирования R и огромную библиотеку R-пакетов машинного обучения и статистики, характеризующихся простотой использования и установки. Эта среда не просто идеальна для изучения R — она будет развиваться вместе с вами, по мере того как вы будете изучать более сложную функциональность языка.

Резюме

Машинное обучение появилось на стыке статистики, управления базами данных и информатики. Это мощный инструмент, способный находить

полезную информацию среди больших объемов данных. Тем не менее, как следует из этой главы, необходимо соблюдать осторожность, чтобы избежать распространенных злоупотреблений ML при решении практических задач.

В основе процесса обучения — абстрагирование данных с преобразованием их в структурированное представление и обобщение структуры с преобразованием в действие, полезность которого затем можно оценить. На практике обучаемая машина использует данные, содержащие примеры и признаки изучаемой концепции, а затем объединяет эти данные в модель, которая используется для прогнозирующих или описательных целей. Эти цели могут быть сгруппированы в задачи, такие как классификация, числовое прогнозирование, обнаружение закономерностей и кластеризация. Среди множества возможных методов нужный алгоритм машинного обучения выбирается в зависимости от входных данных и задачи обучения.

Язык R обеспечивает поддержку ML с помощью пакетов, созданных R-сообществом. Эти мощные инструменты доступны для бесплатной загрузки. Такие пакеты будут представлены в каждой главе книги по мере необходимости.

В следующей главе речь пойдет об основных командах R, которые используются для управления данными и подготовки данных к машинному обучению. У вас может возникнуть желание пропустить эту главу и перейти сразу к практическому применению, однако опыт показывает, что более 80 % времени, затрачиваемого на типичные проекты ML, посвящается этапу подготовки данных, также известному как «выпас данных». Постепенно вы научитесь делать это эффективно, и тогда затраченные усилия окупятся.

2. Управление данными и их интерпретация

Первым этапом любого проекта машинного обучения выступает управление данными и их интерпретация. Возможно, это не так приятно, как создание и развертывание моделей — этапы, на которых вы начинаете видеть плоды своего труда, — однако неразумно игнорировать столь важную подготовительную работу.

Любой обучающий алгоритм хорош ровно настолько, насколько хороши его входные данные, а входные данные зачастую являются сложными, неупорядоченными, разбросанными по нескольким источникам, представленными в разных форматах. Поэтому часто бывает так, что большая часть усилий, затрачиваемых на проект ML, направлена на подготовку и исследование данных.

В этой главе описано три подхода к подготовке данных (data preparation). В первом разделе будут рассмотрены основные структуры данных, которые используются в языке R для хранения данных. Вы близко познакомитесь с этими структурами, когда будете создавать наборы данных и манипулировать ими. Второй раздел является практическим — в нем описано несколько функций, которые используются для ввода и извлечения данных в R. В третьем разделе показаны методы интерпретации данных на примере реального набора данных.

В этой главе рассмотрены следующие темы:

- применение основных структур данных R для хранения и обработки информации;
- простые функции для ввода данных в R из распространенных исходных форматов;
- типичные методы для интерпретации и визуализации сложных данных.

То, как данные представлены в R, будет определять, как с ними нужно работать, поэтому полезно изучить структуры данных R, прежде чем приступать непосредственно к их подготовке. Однако если вы уже знакомы с программированием на языке R, то смело переходите к разделу, посвященному предварительной обработке данных.

Структуры данных R

В разных языках программирования существует множество типов структур данных, каждая из которых имеет свои сильные и слабые стороны и подходит для определенного круга задач. Поскольку R — это язык программирования, широко применяемый для статистического поиска закономерностей, используемые им структуры данных разработаны именно для этого типа задач.

Структуры данных R, наиболее часто встречающиеся в машинном обучении, — это векторы (vectors), факторы (factors), списки (lists), массивы (arrays), матрицы (matrices) и таблицы (фреймы) данных (data frames). Каждая из этих структур адаптирована к конкретной задаче управления данными, поэтому важно понять, как они будут взаимодействовать в вашем R-проекте. В следующих разделах мы рассмотрим сходства этих структур данных и различия между ними.

Векторы

Основной структурой данных в R является *вектор*. В нем хранится упорядоченное множество значений, называемых *элементами*. Вектор может содержать любое количество элементов. Однако все элементы вектора должны быть одного типа; например, вектор не может содержать и цифры, и текст. Тип вектора `v` определяется с помощью команды `typeof(v)`.

В машинном обучении обычно используются следующие типы векторов: `integer` (числа без десятичных знаков), `double` (числа со значениями после запятой), `character` (текстовые данные) и `logical` (принимающие значения `TRUE` или `FALSE`). Есть также два специальных значения: `NA` — отсутствующее (пропущенное или недоступное) значение и `NULL` — нет значения (пустое множество). Два последних значения могут показаться синонимичными, однако на самом деле они различаются. Значение `NA` является заполнителем для чего-либо другого, и поэтому его длина равна единице, а значение `NULL` — действительно пустое, и его длина равна нулю.



Некоторые функции R возвращают векторы типа `integer` и `double` как `numeric`, в то время как другие функции различают эти два типа. В результате все векторы `double` относятся к типу `numeric`, однако не все векторы `numeric` имеют тип `double`.

Вводить большое количество данных вручную утомительно, однако простые векторы можно создавать с помощью функции объединения `c()`. Вектору также можно присвоить имя с помощью оператора `<-`. Это оператор присваивания языка R, используемый практически так же, как оператор `=` во многих других языках программирования.

Для примера построим набор векторов, содержащих данные о трех пациентах. Мы создадим символьный вектор с названием `subject_name` для хранения имен трех пациентов, числовой вектор `temperature` для хранения температуры тела каждого из пациентов (в градусах Фаренгейта) и логический вектор с названием `flu_status` для хранения диагноза каждого пациента (`TRUE`, если у него грипп, и `FALSE` — если нет). Эти три вектора выглядят так, как показано в следующем фрагменте кода:

```
> subject_name <- c("John Doe", "Jane Doe",  
"Steve Graves")> temperature <- c(98.1, 98.6,  
101.4)> flu_status <- c(FALSE, FALSE, TRUE)
```

Значения в R-векторах сохраняют свою последовательность. Поэтому для получения доступа к данным каждого пациента можно использовать его положение в наборе, начиная с 1, и указывать это число в квадратных скобках (`[` и `]`) после имени вектора. Например, чтобы получить значение температуры для пациентки Джейн Доу, второй в списке, достаточно ввести:

```
> temperature[2][1] 98.6
```

В языке R есть множество методов извлечения данных из векторов. Так, с помощью оператора двоеточия можно получить диапазон значений. Например, чтобы получить температуру тела второго и третьего пациентов, введите следующее:

```
> temperature[2:3][1] 98.6 101.4
```

Чтобы исключить элемент, нужно указать его номер со знаком «минус». Так, для исключения данных о температуре второго пациента введите:

```
> temperature[-2][1] 98.1 101.4
```

Иногда полезно указать логический вектор, каждый элемент которого говорит о том, следует ли включать во множество результатов соответствующий элемент. Например, чтобы включить два первых показания температуры, но исключить третье, введите следующее:

```
> temperature[c(TRUE, TRUE, FALSE)][1] 98.1  
98.6
```

Как вскоре будет видно, вектор является основой для многих других структур данных R. Поэтому знание различных векторных операций имеет решающее значение для работы с данными в R.



Как загрузить пример кода

Примеры кода доступны через GitHub по адресу <https://github.com/dataspelunking/MLwR/>. Здесь вы найдете последнюю версию кода на R, а также обсуждение ошибок и общедоступную вики-страницу. Добро пожаловать в сообщество!

Факторы

Как вы, вероятно, помните из главы 1, номинальные признаки представляют собой характеристику с категориями значений. Для хранения номинальных данных можно использовать вектор символов, однако в R есть структура данных, специально предназначенная для этой цели. *Фактор* — это особый случай вектора, который нужен исключительно для представления категориальных или порядковых переменных. В медицинском наборе данных, который мы создаем, можно использовать фактор для указания пола пациента, для этого есть две категории: мужчины и женщины.

Зачем использовать факторы, а не векторы символов? Одним из преимуществ факторов является то, что метки категорий сохраняются только один раз. Вместо сохранения слов **MALE**, **MALE** и **FEMALE** компьютер сохранит числа **1**, **1** и **2**, что позволит сократить объем памяти, необходимый для хранения значений. Кроме того, разные алгоритмы ML по-разному обрабатывают номинальные и числовые функции. Кодирование категориальных

переменных как факторов гарантирует, что R будет обрабатывать категориальные данные надлежащим образом.



Не следует использовать фактор для символьных векторов, которые в действительности не являются категориальными. Так, если вектор хранит в основном уникальные значения, такие как имена или идентификационные коды, сохраните его как символьный вектор.

Для того чтобы создать фактор из символьного вектора, просто примените функцию `factor()`, например:

```
> gender <- factor(c("MALE", "FEMALE",  
"MALE"))> gender[1] MALE FEMALE MALE  
Levels:  
FEMALE MALE
```

Обратите внимание, что при отображении фактора `gender` выводится дополнительная информация о его уровнях. Уровни (`levels`) представляют собой множество возможных категорий, которые может принимать фактор, в данном случае `MALE` и `FEMALE`.

Создавая факторы, можно добавлять в них дополнительные уровни, которые могут быть не представлены в исходных данных. Предположим, мы создали еще один фактор для группы крови, как показано в следующем примере:

```
> blood <- factor(c("O", "AB",  
"A"), levels = c("A", "B", "AB", "O"))>  
blood[1] O AB A  
Levels: A B AB O
```

Определив фактор `blood`, с помощью параметра `levels` мы указали дополнительный вектор, состоящий из четырех возможных групп крови. Несмотря на то что данные включают в себя только группы крови `O`, `AB` и `A`, в факторе `blood`, как показывают результаты, сохраняются все четыре типа. Сохранение дополнительного уровня позволит добавлять пациентов с другой группой крови. Это также гарантирует, что, если бы мы создали таблицу групп крови, мы бы знали, что тип `B` существует, несмотря на его отсутствие в наших исходных данных.

Фактор как структура данных дает возможность включать информацию о последовательности категорий номинальных переменных, что позволяет создавать порядковые признаки. Например, предположим, что у нас есть данные о выраженности симптомов у пациента, закодированные в порядке возрастания степени тяжести — от легкой до умеренной или тяжелой. Мы указываем наличие порядковых данных, предоставляя уровни фактора в желаемой последовательности, перечисляя их в порядке возрастания — от

самого низкого до самого высокого — и присваивая параметру `ordered` значение `TRUE`, как показано ниже:

```
> symptoms <- factor(c("SEVERE", "MILD",  
"MODERATE"), levels = c("MILD",  
"MODERATE", "SEVERE"), ordered =  
TRUE)
```

Полученный в результате фактор `symptoms` теперь включает в себя информацию о требуемом порядке. В отличие от предыдущих факторов уровни этого фактора разделены символами `<` для обозначения последовательности от `MILD` до `SEVERE`:

```
> symptoms[1] SEVERE MILD MODERATE  
Levels:  
MILD < MODERATE < SEVERE
```

Полезным свойством упорядоченных факторов является то, что для них логические тесты работают так, как требуется. Например, можно проверить, насколько выражены симптомы каждого пациента — серьезнее ли они, чем умеренные:

```
> symptoms > "MODERATE"[1] TRUE FALSE FALSE
```

Алгоритмы машинного обучения, способные моделировать упорядоченные данные, ожидают на входе упорядоченные факторы, поэтому обязательно представляйте данные соответствующим образом.

Списки

Список — это структура данных, очень похожая на вектор тем, что она тоже используется для хранения упорядоченного множества элементов. Однако если для вектора требуется, чтобы все его элементы были одинакового типа, то список позволяет собирать разные типы данных. Благодаря такой гибкости списки часто используются для хранения различных типов входных и выходных данных, а также параметров конфигурации для моделей машинного обучения.

Для того чтобы продемонстрировать списки, рассмотрим построенный нами набор данных о трех пациентах, который хранится в шести векторах. Если мы хотим отобразить все данные для первого пациента, то нам нужно ввести пять R-команд:

```
> subject_name[1][1] "John Doe">  
temperature[1][1] 98.1> flu_status[1][1] FALSE>  
gender[1][1] MALE  
Levels: FEMALE MALE> blood[1][1]  
O  
Levels: A B AB O> symptoms[1][1] SEVERE  
Levels:  
MILD < MODERATE < SEVERE
```

Чтобы иметь возможность исследовать данные о пациенте повторно, не вводя заново эти команды, можно использовать список. Он позволяет

сгруппировать все значения в один объект, который можно использовать повторно.

Подобно тому как вектор создается с помощью функции `c()`, список создается с помощью функции `list()`, как показано в следующем примере. Единственное заметное различие заключается в том, что при создании списка каждому компоненту в последовательности должно быть присвоено имя. Имена необязательны, однако их наличие впоследствии позволит получить доступ к значениям по имени, а не по порядковому номеру. Для того чтобы создать список с именованными компонентами для всех данных первого пациента, введите следующую команду:

```
> subject1 <- list(fullname =
subject_name[1], temperature =
temperature[1], flu_status =
flu_status[1], gender =
gender[1], blood =
blood[1], symptoms =
symptoms[1])
```

Теперь данные пациента собраны в списке `subject1`:

```
> subject1$fullname[1] "John
Doe"$temperature[1] 98.1$flu_status[1]
FALSE$gender[1] MALELevels: FEMALE MALE$blood[1]
OLevels: A B AB O$symptoms[1] SEVERELevels: MILD
< MODERATE < SEVERE
```

Обратите внимание, что значения помечены именами, которые мы указали в предыдущей команде. Поскольку список сохраняет последовательность значений так же, как и вектор, к его компонентам можно получить доступ по порядковым номерам, как показано далее для значения `temperature`:

```
> subject1[2]$temperature[1] 98.1
```

Результатом применения векторных операторов к объекту списка является другой объект списка, который представляет собой подмножество исходного списка. Например, предыдущий код возвращает список, состоящий из единственного компонента `temperature`. Чтобы вернуть один элемент списка с его **собственным** типом данных, при выборе компонента списка используйте двойные скобки (`[[` и `]]`). Например, следующая команда возвращает числовой вектор единичной длины:

```
> subject1[[2]][1] 98.1
```

Для точности лучше обращаться к компонентам списка по имени, ставя после имени списка знак `$` и имя компонента следующим образом:

```
> subject1$temperature[1] 98.1
```

Подобно записи в двойных скобках, эта запись возвращает компонент списка его собственного типа данных (в данном случае — числовой вектор единичной длины).



Доступ к значению по имени также гарантирует, что всегда будет получен правильный элемент, даже если впоследствии порядок элементов списка изменится.

Для того чтобы получить несколько элементов списка, нужно указать вектор имен. Следующая команда возвращает подмножество списка `subject1`, который содержит только компоненты `temperature` и `flu_status`:

```
> subject1[c("temperature",  
"flu_status")]$temperature[1] 98.1$flu_status[1]  
FALSE
```

С помощью списков и списков списков можно строить целые наборы данных. Например, можно создать списки `subject2` и `subject3` и сгруппировать их в списочный объект с именем `pt_data`. Такое построение набора данных широко распространено, но в языке R есть специальная структура данных.

Фреймы данных

Наиболее важной структурой данных R, используемой в ML, является, безусловно, *фрейм данных* — структура, аналогичная электронной таблице или базе данных, в которой есть строки и столбцы. В контексте языка R фрейм данных можно представить как список векторов или факторов, каждый из которых имеет одинаковое количество значений. Поскольку фрейм данных является, в сущности, списком объектов векторного типа, он объединяет в себе свойства как векторов, так и списков.

Создадим фрейм данных для нашего набора данных о пациентах. Для этого используем созданные ранее векторы данных о пациенте и вызовем функцию `data.frame()`, которая объединит их во фрейм данных:

```
> pt_data <- data.frame(subject_name,  
temperature, flu_status,  
gender, blood,  
symptoms, stringsAsFactors  
= FALSE)
```

Как вы могли заметить, в этом коде появилось нечто новое. Мы включили сюда дополнительный параметр: `stringsAsFactors=FALSE`. Если его не указать, то R автоматически преобразует каждый символьный вектор в фактор.

Иногда это свойство бывает полезным, но его применение не всегда оправданно. В данном случае, например, поле `subject_name` определенно не считается категориальными данными, поскольку имена не являются категориями значений. Поэтому, присвоив параметру `stringsAsFactors` значение `FALSE`, мы будем преобразовывать символьные векторы в факторы только в тех случаях, когда это имеет смысл для проекта.

При выводе фрейма данных `pt_data` мы видим, что его структура сильно отличается от структур данных, с которыми мы уже имели дело:

```
>
pt_data  subject_name  temperature  flu_status
gender  blood  symptoms1  John
Doe      98.1      FALSE      MALE      O
SEVERE2  Jane
Doe      98.6      FALSE      FEMALE    AB
MILD3   Steve
Graves   101.4      TRUE      MALE      A
MODERATE
```

В отличие от одномерных векторов, факторов и списков фрейм данных имеет два измерения и отображается в матричном формате. Конкретно в этом фрейме данных есть один столбец для каждого вектора данных пациента и одна строка для каждого пациента. В машинном обучении столбцы фрейма данных — это признаки, или атрибуты, а строки — примеры.

Для того чтобы извлечь целые столбцы (векторы) данных, можно воспользоваться тем, что фрейм данных представляет собой просто список векторов. Подобно спискам, самый прямой способ извлечь отдельный элемент — обратиться к нему по имени. Например, чтобы получить вектор `subject_name`, введите следующее:

```
> pt_data$subject_name[1] "John Doe"      "Jane
Doe"      "Steve Graves"
```

Подобно спискам, вектор имен можно использовать для извлечения нескольких столбцов из фрейма данных:

```
> pt_data[c("temperature",
"flu_status")]  temperature  flu_status1
98.1      FALSE2      98.6      FALSE3
101.4      TRUE
```

Если запрашивать столбцы во фрейме данных по имени, то результатом будет фрейм, содержащий все строки данных для указанных столбцов. Команда `pt_data[2:3]` тоже извлечет

столбцы `temperature` и `flu_status`, однако обращение к столбцам по имени приводит к ясному и простому в сопровождении R-коду, который не перестанет работать, если фрейм данных будет впоследствии переупорядочен.

Для извлечения из фрейма данных конкретных значений используются методы, подобные тем, что применяются для доступа к значениям векторов. Однако здесь есть важное отличие: поскольку фрейм данных является двумерным, необходимо указывать не только нужные строки, но и столбцы. Сначала указываются строки, затем ставится запятая, за которой следуют столбцы: `[row, col]`. Как и для векторов, номера строк и столбцов начинаются с единицы.

Например, чтобы извлечь значение, стоящее в первой строке и втором столбце фрейма данных пациента, используется следующая команда:

```
> pt_data[1, 2][1] 98.1
```

Если требуется извлечь несколько строк или столбцов данных, то нужно указать векторы, соответствующие желаемым строкам и столбцам.

Следующий оператор извлекает данные из первой и третьей строк, а также из второго и четвертого столбцов:

```
> pt_data[c(1, 3), c(2,
4)] temperature gender1 98.1 MALE
3 101.4 MALE
```

Чтобы сослаться на каждую строку или столбец, просто оставьте часть строки или столбца пустой. Например, чтобы извлечь все строки первого столбца, введите команду:

```
> pt_data[, 1][1] "John Doe" "Jane
Doe" "Steve Graves"
```

Так можно извлечь все столбцы для первой строки:

```
> pt_data[1,
] subject_name temperature flu_status g
ender blood symptoms1 John
Doe 98.1 FALSE MALE
0 SEVERE
```

А так — все данные:

```
> pt_data[ ,
] subject_name temperature flu_status gend
er blood symptoms1 John
Doe 98.1 FALSE MALE 0
SEVERE2 Jane
Doe 98.6 FALSE FEMALE AB
MILD3 Steve
```

```
Graves          101.4          TRUE          MALE
A      MODERATE
```

Конечно, к столбцам лучше обращаться по имени, а не по порядковому номеру, а для исключения строк или столбцов данных можно использовать отрицательные значения. Таким образом, результат команды:

```
> pt_data[c(1, 3), c("temperature",
"gender")] temperature gender1          98.1
MALE3          101.4          MALE
```

эквивалентен результату следующей:

```
> pt_data[-2, c(-1, -3, -5, -
6)] temperature gender1          98.1          MALE
3          101.4          MALE
```

Иногда во фреймах данных необходимо создавать новые столбцы — например, как функцию от существующих столбцов. Так, нам может потребоваться преобразовать показания температуры во фрейме данных о пациенте из шкалы Фаренгейта в шкалу Цельсия. Для этого достаточно воспользоваться оператором присваивания и присвоить результат преобразования столбцу с новым именем следующим образом:

```
> pt_data$temp_c <- (pt_data$temperature - 32)
* (5 / 9)
```

Чтобы убедиться в правильности вычислений, сравним новый столбец `temp_c` с температурой по Цельсию с предыдущим столбцом `temperature`, где указана температура по Фаренгейту:

```
> pt_data[c("temperature",
"temp_c")] temperature          temp_c1          98.1
36.722222          98.6          37.000003          101.4
38.55556
```

Увидев эти показания рядом, мы можем подтвердить, что расчет выполнен верно.

Чтобы лучше ознакомиться с фреймами данных, попробуйте выполнить аналогичные операции с набором данных о пациентах или, что еще лучше, использовать данные одного из ваших собственных проектов. Важно понимать суть этих операций, чтобы выполнить большую часть работы, предстоящей в следующих главах.

Матрицы и массивы

Кроме фреймов данных, в R есть и другие структуры, которые хранят значения в табличной форме. *Матрица* — это структура данных, сведенных в двумерную таблицу, где данные представлены в виде строк и столбцов. Подобно векторам, матрицы в R могут содержать только один

тип данных, хотя они чаще всего используются для математических операций, и поэтому обычно хранят только числа.

Для того чтобы создать матрицу, достаточно вызвать функцию `matrix()` и указать вектор данных, а также параметр, определяющий количество строк (`nrow`) или столбцов (`ncol`). Например, для создания матрицы 2×2 , в которой хранятся числа от единицы до четырех, можно использовать параметр `nrow`, чтобы запросить данные, которые будут разделены на две строки:

```
> m <- matrix(c(1, 2, 3, 4), nrow = 2)>
m      [,1] [,2][1,] 1 3[2,] 2 4
```

Или же можно создать такую же матрицу с использованием параметра `ncol=2`:

```
> m <- matrix(c(1, 2, 3, 4), ncol = 2)>
m      [,1] [,2][1,] 1 3[2,] 2 4
```

Вероятно, вы заметили, что в R сначала загружается первый столбец матрицы, а затем — второй. Это называется *размещением по столбцам*, которое применяется в R по умолчанию для загрузки матриц.



Чтобы изменить способ размещения, используемый по умолчанию, и загрузить матрицу по строкам, при создании матрицы установите параметр `byrow = TRUE`.

Рассмотрим на примере, что произойдет, если увеличить количество значений матрицы.

Если указать шесть значений, то для двух строк будет создана матрица из трех столбцов:

```
> m <- matrix(c(1, 2, 3, 4, 5, 6), nrow = 2)>
m      [,1] [,2] [,3][1,] 1 3 5[2,] 2 4 6
```

Если задать два столбца, то будет создана матрица из трех строк:

```
> m <- matrix(c(1, 2, 3, 4, 5, 6), ncol = 2)>
m      [,1] [,2][1,] 1 4[2,] 2 5[3,] 3 6
```

Как и в случае с фреймами данных, значения матриц могут быть извлечены с использованием нотации `[row, column]`.

Например, `m[1, 1]` для матрицы `m` вернет значение `1`, а `m[3, 2]` — `6`.

Кроме того, можно извлекать целые строки или столбцы:

```
> m[1, ][1] 1 4> m[, 1][1] 1 2 3
```

С матричной структурой тесно связаны *массивы* — многомерные таблицы данных. Если матрица состоит из строк и столбцов значений, то в массиве есть строки, столбцы и множество дополнительных слоев

значений. В последующих главах матрицы иногда будут встречаться, а вот использовать массивы в рамках этой книги нет необходимости.

Управление данными в R

Одна из проблем, с которыми приходится сталкиваться при работе со сверхбольшими наборами данных, заключается в сборе, подготовке данных, поступающих из различных источников, и управлении ими. Далее на примере реальных задач машинного обучения мы подробно рассмотрим подготовку данных, их очистку и управление ими. В этом разделе описаны лишь основные функции для ввода и вывода данных в R.

Сохранение, загрузка и удаление структур данных в R

Если вы один раз потратили много времени на преобразование фрейма данных в нужную форму, то вам не нужно будет заново все переделывать при каждом перезапуске R-сеанса. Чтобы сохранить структуру данных в файле, который впоследствии можно загрузить или перенести в другую систему, используйте функцию `save()`. Функция записывает одну или несколько структур данных R в местоположение, указанное параметром `file`. Файлы данных R имеют расширение `.RData`.

Предположим, у вас есть три объекта с именами `x`, `y` и `z`, которые вы хотели бы сохранить в постоянном файле. Независимо от того, являются ли эти объекты векторами, факторами, списками или фреймами данных, их можно сохранить в файле `mydata.RData` с помощью следующей команды:

```
> save(x, y, z, file = "mydata.RData")
```

Команда `load()` позволяет воссоздать любые структуры данных, которые были сохранены в файле `.RData`. Для того чтобы загрузить созданный файл `mydata.RData`, просто введите следующее:

```
> load("mydata.RData")
```

По этой команде в среде R будут воссозданы структуры данных `x`, `y` и `z`.



Будьте осторожны, загружая данные! Все структуры данных, хранящиеся в файле, который импортируется с помощью команды `load()`, будут добавлены в ваше рабочее пространство, даже если они перезаписывают что-то еще, над чем вы работаете.

Если нужно быстро завершить сеанс R, воспользуйтесь командой `save.image()` — она запишет весь сеанс в файл, называемый `.RData`. По умолчанию при следующем запуске R будет

искать этот файл, и сеанс будет воссоздан точно в том виде, в каком вы его оставили.

После того как вы некоторое время поработаете в сеансе R, вы, возможно, накопите несколько структур данных.

Функция `ls()` возвращает вектор, который состоит из всех структур данных, находящихся в памяти в настоящий момент. Например, если вы выполняли все команды этой главы, то функция `ls()` возвратит следующее:

```
>
ls()[1] "blood" "flu_status" "gender"
      "m" [5] "subject_name" "subject1" "symp
toms" [9] "temperature"
```

При выходе из сеанса R автоматически удаляет из памяти все структуры данных, однако для больших объектов, возможно, вы захотите освободить память раньше. Для этого можно использовать функцию удаления `rm()`. Например, чтобы удалить объекты `m` и `subject1`, просто введите следующее:

```
> rm(m, subject1)
```

Функции `rm()` также можно передать символьный вектор имен объектов, подлежащих удалению. Если же передать ей результат работы функции `ls()`, то будет очищен весь сеанс R:

```
> rm(list = ls())
```

Будьте очень осторожны при выполнении этой команды, так как перед удалением объектов вы не получите предупреждения!

Импорт и сохранение данных из CSV-файлов

Общедоступные наборы данных обычно хранятся в текстовых файлах. Текстовые файлы могут быть прочитаны практически на любом компьютере и в любой операционной системе, что делает этот формат универсальным. Текстовые файлы также могут быть экспортированы и импортированы в такие программы, как Microsoft Excel, что позволяет быстро и легко работать с данными, представленными в виде электронных таблиц.

Табличный (представленный в виде таблицы) файл данных имеет матричную структуру: каждая строка текста отражает один пример и каждый пример имеет одинаковое количество признаков. Значения признаков в каждой строке разделены предопределенным символом, который называется *разделителем*. В первой строке файла табличных данных часто содержатся имена столбцов. Эта строка называется *заголовком*.

Пожалуй, самым распространенным форматом табличного текстового файла является файл *значений, разделенных запятыми* (comma-separated values, CSV), в котором, как следует из названия, в качестве разделителя используются запяты. CSV-файлы можно импортировать и экспортировать во многих распространенных приложениях. CSV-файл, описывающий созданный ранее набор медицинских данных, может быть сохранен следующим образом:

```
subject_name, temperature, flu_status, gender, blood_type
John Doe, 98.1, FALSE, MALE, O
Jane Doe, 98.6, FALSE, FEMALE, AB
Steve Graves, 101.4, TRUE, MALE, A
```

Пусть файл с данными о пациентах `pt_data.csv` расположен в рабочем каталоге R. Тогда функцию `read.csv()` можно использовать следующим образом:

```
> pt_data <- read.csv("pt_data.csv",
stringsAsFactors = FALSE)
```

Эта функция прочитает CSV-файл и поместит данные во фрейм с именем `pt_data`. Подобно тому как мы сделали это при создании фрейма данных, здесь нужно использовать параметр `stringsAsFactors=FALSE`, чтобы предотвратить преобразование всех текстовых переменных в факторы. Если вы уверены, что все столбцы в CSV-файле действительно являются факторами, данный шаг лучше выполнить самостоятельно, а не предоставлять это R.



Если набор данных находится за пределами рабочего каталога R, то при вызове функции `read.csv()` нужно указать полный путь к CSV-файлу (например, `"/path/to/mydata.csv"`).

По умолчанию R предполагает, что в CSV-файле есть строка заголовка, в которой перечислены имена объектов из набора данных. Если в CSV-файле нет заголовка, укажите параметр `header=FALSE`, как показано в следующей команде, и R назначит имена признаков по умолчанию, нумеруя их `V1`, `V2` и т.д.:

```
> mydata <- read.csv("mydata.csv",
stringsAsFactors =
FALSE,
header = FALSE)
```

Функция `read.csv()` является частным случаем функции `read.table()`, которая позволяет считывать табличные данные в различных форматах, включая другие форматы с разделителями, такие как *значения, разделенные табуляцией* (tab-separated values, TSV).

Подробнее о семействе функций `read.table()` читайте на странице справочной системы R, открыв ее с помощью команды `?read.table`.

Чтобы сохранить фрейм данных в CSV-файле, используйте функцию `write.csv()`. Если фрейм данных называется `pt_data`, просто введите следующее:

```
> write.csv(pt_data, file = "pt_data.csv",  
row.names = FALSE)
```

По этой команде в рабочий каталог R будет записан CSV-файл с именем `pt_data.csv`. Параметр `row.names` переопределяет заданное по умолчанию поведение R, которое заключается в выводе имен строк в CSV-файле. Как правило, эти выходные данные не нужны и будут лишь увеличивать размер получаемого файла.

Исследование данных и их интерпретация

После сбора данных и загрузки их в структуры данных R следующий шаг в процессе машинного обучения — подробное исследование данных. Именно на этом этапе начинается исследование признаков и примеров данных, а также выявление особенностей, которые делают эти данные уникальными. Чем лучше будут интерпретированы данные, тем лучше вы сможете выбрать модель ML для вашей задачи.

Лучший способ изучить процесс исследования данных — рассмотреть его на примере. В этом разделе рассмотрен набор данных `usedcars.csv`, который содержит реальные данные о подержанных автомобилях, выставившихся на продажу на популярном американском сайте в 2012 году.



Набор данных `usedcars.csv` можно скачать вместе с файлами примеров к книге с GitHub. Если вы решили выполнить этот пример, сначала загрузите и сохраните указанный файл в рабочем каталоге R.

Поскольку набор данных хранится в форме CSV, мы можем использовать функцию `read.csv()` для загрузки данных из него во фрейм данных R:

```
> usedcars <- read.csv("usedcars.csv",  
stringsAsFactors = FALSE)
```

Теперь, когда у нас есть фрейм данных `usedcars`, представьте себя в роли специалиста по анализу данных, который должен изучить данные о подержанных автомобилях. Несмотря на то что изучение данных — гибкий процесс, его этапы можно представить как своего рода исследование, в котором даются ответы на вопросы о данных. Вопросы могут различаться в зависимости от проекта, однако типы вопросов всегда похожи. Вы

сможете адаптировать основные этапы этого исследования к любому набору данных, будь он большим или маленьким.

Структуры данных

Первые вопросы, которые следует задать при исследовании нового набора данных, должны быть о том, как организован набор данных. Если повезет, то ваш источник предоставит *словарь данных* — документ, в котором описаны признаки набора данных. В нашем случае данные о подержанных автомобилях не сопровождаются таким документом, поэтому придется создать его самостоятельно.

Функция `str()` представляет собой метод отображения структуры R-объектов, таких как фреймы данных, векторы и списки. Эту функцию можно использовать для создания базовой схемы словаря данных:

```
> str(usedcars) 'data.frame' :      150 obs. of
6 variables:$ year      : int  2011 2011
2011 2011 ...$ model    : chr  "SEL"
"SEL" "SEL" "SEL" ...$ price      :
int  21992 20995 19995 17809 ...$
mileage      : int  7413 10926 7351 11613
...$ color   : chr  "Yellow" "Gray"
"Silver" "Gray" ...$ transmission :
chr  "AUTO" "AUTO" "AUTO" "AUTO" ...
```

Как много информации о наборе данных мы получили с помощью такой простой команды! Фрагмент `150obs` означает, что данные включают в себя 150 наблюдений — еще один способ сообщить, что набор данных содержит 150 записей, или примеров. Количество наблюдений часто обозначается буквой n . Поскольку мы знаем, что данные описывают подержанные автомобили, то можем теперь предположить, что у нас есть примеры $n = 150$ автомобилей для продажи.

Фрагмент `6variables` говорит о наличии шести признаков, которые записаны в данных. Эти признаки перечислены по названиям в отдельных строках. Глядя на строку для признака `color`, можно отметить дополнительные детали:

```
$ color      : chr  "Yellow" "Gray"
"Silver" "Gray" ...
```

После имени переменной стоит метка `chr`, которая сообщает, что этот признак является символьным. В этом наборе данных три переменные символьные, а три отмечены как `int`, что означает тип `integer`. Набор данных `usedcars` содержит только символьные и целочисленные переменные, однако в других случаях, при использовании не только

целочисленных данных, можно встретить тип `num` или `numeric`. Все факторы будут иметь тип `factor`. После типа каждой R-переменной указана последовательность первых нескольких значений признаков. Первыми четырьмя значениями признака `color` являются `"Yellow"`, `"Gray"`, `"Silver"` и `"Gray"`.

Применив к именам и значениям признаков некоторые знания из предметной области, можно сделать предположения о том, что представляют собой переменные. Переменная `year` может указывать на год выпуска автомобиля или на год, когда было опубликовано объявление о продаже. Позже нужно будет исследовать этот признак более подробно, поскольку четыре примера значений (`2011`, `2011`, `2011`, `2011`) могут послужить аргументом в пользу любой из этих версий.

Переменные `model`, `price`, `mileage`, `color` и `transmission`, скорее всего, относятся к характеристикам продаваемого автомобиля.

Похоже, что данным были присвоены осмысленные имена переменных, однако это не всегда так. Иногда наборы данных имеют признаки с бессмысленными именами или кодами, такими как `V1`. В этих случаях может потребоваться дополнительное расследование, чтобы определить, что представляет собой тот или иной признак. Но даже с информативными именами признаков следует всегда скептически относиться к предоставленным меткам. Продолжим исследование.

Числовые переменные

Чтобы исследовать числовые переменные в данных о подержанных автомобилях, мы воспользуемся базовым набором измерений, описывающим значения, известные как *сводная статистика*.

Функция `summary()` позволяет вывести несколько видов сводной статистики. Рассмотрим один из признаков — `year`:

```
> summary(usedcars$year)Min. 1st
Qu.  Median  Mean  3rd
Qu.  Max.2000  2008  2009  2009
2010  2012
```

Пока не будем обращать внимание на конкретные значения — важен тот факт, что здесь есть цифры `2000`, `2008` и `2009`, которые говорят о том, что переменная `year` указывает на год выпуска автомобиля, а не на год, когда было опубликовано объявление, поскольку, как известно, списки автомобилей были составлены в 2012 году.

Если передать функции `summary()` вектор имен столбцов, то получим сводную статистику сразу для нескольких числовых переменных:

```

> summary(usedcars[c("price",
"mileage")])
      price      mileageMin.      :
3800   Min.      : 48671st Qu. :10995   1st Qu. :
27200Median      :13592      Median      :
36385Mean        :12962      Mean        : 442613rd Qu.
:14904      3rd Qu. :
55125Max.        :21992      Max.        :151479

```

Сводная статистика, предоставленная функцией `summary()`, — простой, но мощный инструмент для исследования данных. Она бывает двух типов: меры центра и меры разброса.

Измерение средних значений: среднее арифметическое и медиана

Измерение *средних значений* — это задача из статистики, используемая для определения значения, которое попадает в середину набора данных. Скорее всего, один распространенный показатель центра вам уже знаком: среднее значение. В целом, когда что-то считается средним, это значит, что оно находится где-то между крайними точками шкалы.

Если у ученика средний балл успеваемости, это говорит о том, что ему были выставлены средние отметки в классе. Средний — не слишком легкий, не слишком тяжелый. В целом средний предмет типичен и не особо отличается от других в своей группе. Его можно представить как образец, по которому судят обо всех остальных.

В статистике среднее также называется *средним арифметическим* — это измерение, определяемое как сумма всех значений, разделенная на количество значений. Например, для того, чтобы вычислить средний доход в группе из трех человек с доходами 36 000, 44 000 и 56 000 долларов, можно выполнить следующую команду:

```

> (36000 + 44000 + 56000) / 3[1] 45333.33

```

В R также есть функция `mean()`, которая вычисляет среднее арифметическое вектора чисел:

```

> mean(c(36000, 44000, 56000))[1] 45333.33

```

Средний доход этой группы людей составляет примерно 45 333 доллара. Концептуально это можно представить как доход, который был бы у каждого человека, если бы общая сумма дохода делилась поровну на всех.

Напомню, что при выводе предыдущих результатов `summary()` были заданы средние значения переменных `price` и `mileage`. Эти значения указывают на то, что типичная подержанная машина в этом наборе данных стоит 12 962 доллара, а показание ее одометра — 44 261. Что это говорит о данных? Поскольку средняя цена относительно низкая, то можно ожидать, что в наборе данных содержатся автомобили экономкласса. Конечно, среди них могут также присутствовать автомобили повышенной комфортности

последней модели с большим пробегом. Но сравнительно низкий статистический показатель среднего пробега не приводит доказательств в поддержку этой гипотезы. Впрочем, он также не предоставляет доказательств, позволяющих проигнорировать такую возможность. Следует помнить об этом при дальнейшем изучении данных.

Среднее арифметическое является наиболее часто упоминаемым статистическим показателем для измерения среднего значения в наборе данных, однако не всегда самым подходящим. Другим часто используемым показателем среднего значения выступает *медиана* — значение, которое находится в середине упорядоченного списка значений. Как и в случае со средним арифметическим, в R есть функция `median()`, которую можно применить к нашим данным о зарплате следующим образом:

```
> median(c(36000, 44000, 56000))[1] 44000
```

Итак, поскольку среднее значение равно `44000`, то медианный доход составляет 44 000 долларов.



Если в наборе данных четное число значений, то среднего значения не существует. В этом случае медиана обычно рассчитывается как среднее из двух значений, находящихся в центре упорядоченного списка. Например, медиана значений 1, 2, 3 и 4 составляет 2,5.

На первый взгляд кажется, что медиана и среднее арифметическое очень похожи. Конечно, среднее арифметическое 45 333 доллара и медиана 44 000 долларов не очень различаются. Зачем же нужны два измерения средних значений? Причина заключается в том, что среднее арифметическое и медиана по-разному зависят от значений, находящихся на границах диапазона. В частности, среднее арифметическое очень чувствительно к *выбросам* — значениям, которые нетипично велики или малы по сравнению с большинством данных. Таким образом, поскольку среднее арифметическое чувствительнее к выбросам, оно с большей вероятностью будет сдвигаться в большую или меньшую сторону при наличии небольшого числа экстремальных значений.

Снова обратимся к полученным медианным значениям в выводе результатов `summary()` для набора данных о подержанных автомобилях. Среднее арифметическое и медиана для цены довольно близки (различаются примерно на 5 %), однако разница между средним арифметическим и медианой пробега намного больше. Среднее арифметическое пробега составляет 44 261, что почти на 20 % больше медианы — 36 385. Поскольку среднее арифметическое более чувствительно к экстремальным значениям, чем медиана, тот факт, что среднее арифметическое получилось намного больше медианы, может

навести на мысль, что в наборе данных есть подержанные автомобили с чрезвычайно высокими значениями пробега. Чтобы подробнее исследовать этот факт, нужно добавить в анализ дополнительную сводную статистику.

Измерение разброса: квартили и пятичисловая сводка

Среднее арифметическое и медиана позволяют быстро получить сводные значения, однако эти измерения средних значений мало скажут о том, существует ли разнообразие в измерениях. Для того чтобы измерить разнообразие, нужно использовать другой тип сводной статистики, касающийся *разброса* данных, — того, насколько тесно или широко распределены значения. Знание о разбросе дает представление о максимумах и минимумах данных, а также о том, насколько большинство значений близко к среднему арифметическому и медиане.

Пятичисловая сводка — это набор из пяти статистических показателей, которые приблизительно отображают разброс значений некоего признака. Все эти пять статистических показателей включены в выходные данные функции `summary()`. Перечислим их по порядку:

- минимум (`Min.`);
- первый квартиль, или Q1 (`1stQu.`);
- медиана, или Q2 (`Median`);
- третий квартиль, или Q3 (`3rdQu.`);
- максимум (`Max.`).

Как и следовало ожидать, минимальное и максимальное значения являются наиболее экстремальными значениями признака, указывая на наименьшее и наибольшее значения соответственно. Для вычисления этих значений для вектора в R есть функции `min()` и `max()`.

Интервал между минимумом и максимумом называется *диапазоном*. В R есть функция `range()`, которая возвращает минимальное и максимальное значение:

```
> range(usedcars$price)[1] 3800 21992
```

Сочетание `range()` и разностной функции `diff()` позволяет вычислить статистику диапазона в одной строке кода:

```
> diff(range(usedcars$price))[1] 18192
```

Первый и третий квартили, Q1 и Q3, — это значения, ниже и выше которых расположена четверть всех значений. Наряду с медианой (Q2) *квартили* делят набор данных на четыре части, на каждую из которых приходится одинаковое количество значений.

Квартили являются частным случаем *квантилей*, представляющих собой числа, которые делят данные на одинаковые по размеру множества. Кроме квартилей, в число часто используемых квантилей

входят *тертили* (деление на три части), *квинтили* (пять частей), *децили* (десять частей) и *процентили* (сто частей).



Процентили часто используются для описания ранжирования значений. Например, студент, тестовый балл которого оценен на 99-й процентиль, показал результаты, которые как минимум не хуже, чем у 99 % других тестируемых.

Средние 50 % данных между первым и третьим квартилями представляют особый интерес, потому что это простая мера разброса. Разница между Q1 и Q3 называется *межквартильным диапазоном* (interquartile range, IQR) и может быть вычислена с помощью функции `IQR()`:

```
> IQR(usedcars$price)[1] 3909.5
```

Мы могли бы вычислить это значение для переменной `usedcars$price` вручную, по результатам `summary()`: $14\ 904 - 10\ 995 = 3909$. Небольшая разница между нашими расчетами и результатом `IQR()` обусловлена тем, что R при выводе результатов `summary()` автоматически округляет числа.

Функция `quantile()` представляет собой универсальный инструмент для определения квантилей в наборе значений. По умолчанию функция `quantile()` возвращает пятичисловую сводку. Если применить эту функцию к переменной `usedcars$price`, то получим ту же сводную статистику, что и раньше:

```
>
quantile(usedcars$price)      0%      25%
50%      75%      100% 3800.0    10995.0    135
91.5    14904.5    21992.0
```



Для вычисления квантилей предусмотрено много методов обработки взаимосвязей между наборами значений, не имеющих единого среднего значения. Функция `quantile()` позволяет выбрать один из девяти алгоритмов разрешения конфликтов, указав параметр `type`. Если в проекте требуется точно определить квантиль, важно изучить документацию к этой функции с помощью команды `?quantile`.

Указав дополнительный параметр `probs` и используя вектор, содержащий точки деления, можно получить произвольные квантили — например, 1-й и 99-й процентили:

```
> quantile(usedcars$price, probs = c(0.01,
0.99))      1%      99% 5428.69    20505.00
```

Функция последовательности `seq()` генерирует векторы с равномерно распределенными значениями. Это облегчает получение других фрагментов данных, таких как квинтили (пять групп), показанные в следующем примере:

```
> quantile(usedcars$price, seq(from = 0, to =  
1, by =  
0.20))
```

	0%	20%	40%	60%
		3800.0	10759.4	12993.8
	13992.0	14999.0	21992.0	

Теперь, зная, что такое пятичисловая сводка, мы можем пересмотреть сводную статистику о подержанных машинах, полученную с помощью `summary()`. Что касается переменной `price`, то ее минимальное значение составляет 3800 долларов, а максимальное — 21992. Интересно, что разница между минимумом и Q1 составляет около 7000 долларов, как и разница между Q3 и максимумом; однако разница между Q1 и медианой до Q3 составляет примерно 2000 долларов. Это говорит о том, что нижние и верхние 25 % значений разбросаны более широко, чем средние 50 % значений, которые, по-видимому, более тесно сгруппированы вокруг центра. Аналогичная тенденция наблюдается для переменной `mileage`. Как вы узнаете позже из этой главы, данный паттерн разброса достаточно широко распространен, и его называют нормальным распределением данных.

Разброс переменной `mileage` выявляет еще одно интересное свойство — разница между Q3 и максимумом намного больше, чем между минимумом и Q1. Другими словами, разброс между большими значениями пробега гораздо больше, чем между малыми.

Это открытие помогает объяснить, почему среднее арифметическое оказалось намного больше, чем медиана. Поскольку среднее арифметическое чувствительно к экстремальным значениям, оно оказалось выше, в то время как медиана осталась на том же уровне. Это важное свойство, которое становится более очевидным, когда данные представлены визуально.

Визуализация числовых переменных: диаграммы размаха

Визуализация числовых переменных может быть полезна для выявления проблем с данными. Распространенным способом визуализации пятичисловой сводки является *диаграмма размаха*, также известная как «ящик с усами» (*box-and-whisker*). На диаграмме размаха отображаются центр и разброс числовой переменной в виде, позволяющем быстро получить представление о диапазоне и отклонении переменной или же сравнить ее с другими переменными.

Рассмотрим диаграмму размаха для данных о цене и пробеге подержанных автомобилей. Чтобы получить диаграмму размаха для переменной, воспользуемся функцией `boxplot()`. Зададим пару дополнительных параметров, `main` и `ylab`, чтобы добавить к диаграмме заголовков и обозначить ось Y (вертикальную ось) соответственно. Команды построения диаграмм размаха для переменных `price` и `mileage` выглядят так:

```
> boxplot(usedcars$price, main = "Boxplot of  
Used Car Prices", ylab = "Price ($)")>  
boxplot(usedcars$mileage, main = "Boxplot of Used  
Car Mileage", ylab = "Odometer (mi.)")
```

R выдаст следующие изображения (рис. 2.1).

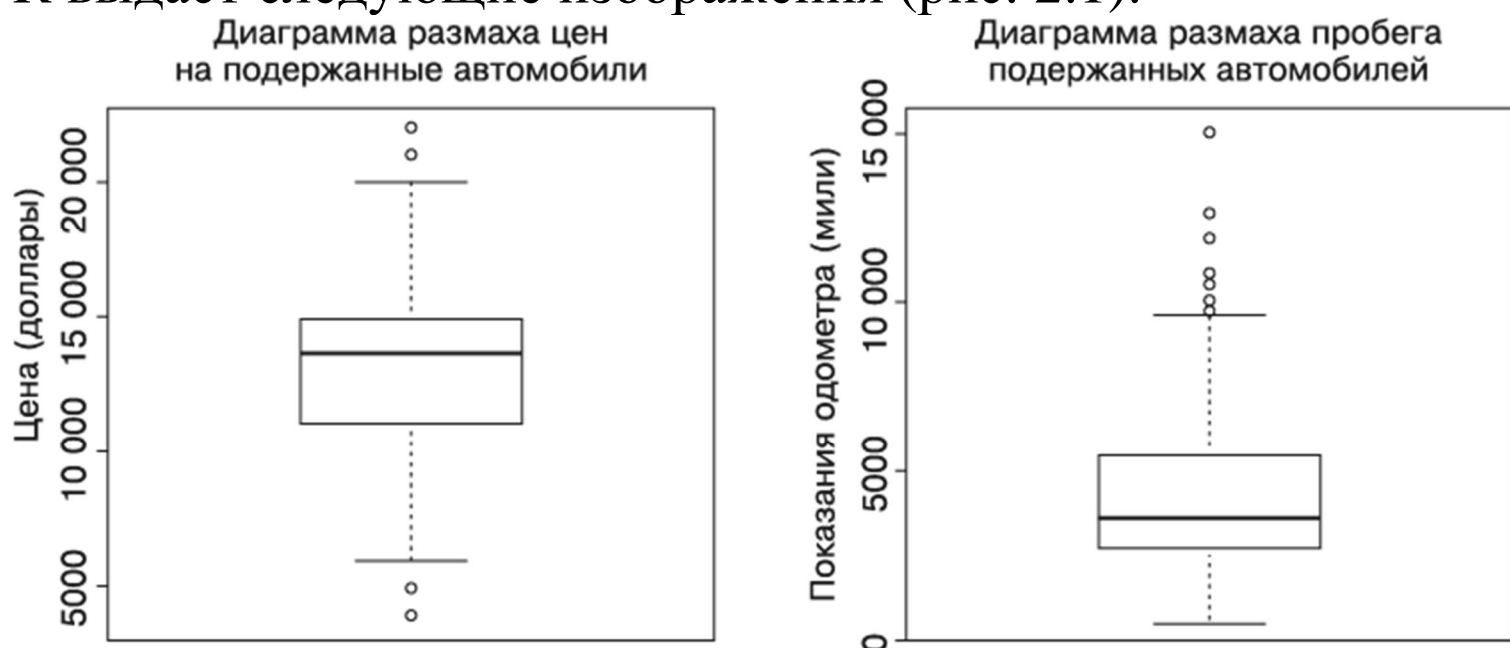


Рис. 2.1. Диаграммы размаха для данных о цене и пробеге подержанных автомобилей

Диаграмма размаха отображает пятичисловую сводку, используя горизонтальные линии и точки. Горизонтальные линии, образующие прямоугольник в середине каждого рисунка, отражают значения Q1, Q2 (медиана) и Q3, если смотреть на диаграмму снизу вверх. Медиана обозначена темной линией, которая соответствует значениям 13 592 доллара по вертикальной оси для цены и 36 385 миль для пробега.



На простых диаграммах размаха, таких как диаграммы, показанные на рис. 2.1, ширина «ящика» является произвольной и не отображает какую-либо характеристику данных. Для более сложных видов анализа можно варьировать форму и размер «ящиков», чтобы облегчить сравнение данных по нескольким группам. Если вы хотите узнать больше о таких функциях, начните с изучения параметров `notch` и `varwidth` в документации к R-функции `boxplot()`, введя команду `?boxplot`.

Минимальные и максимальные значения можно проиллюстрировать с помощью «усов», которые простираются ниже и выше «ящика»; однако распространенное соглашение допускает длину «усов» не более чем на 1,5

IQR ниже Q1 или выше Q3. Любые значения, которые выходят за пределы этого порога, считаются выбросами и обозначаются кружками либо точками. Так, напомним, что IQR для переменной `price` составляет 3909, притом что Q1 равно 10 995 и Q3 равно 14 904. Таким образом, выбросом является любое значение, меньшее чем $10\,995 - 1,5 \times 3909 = 5131,5$ или большее чем $14\,904 + 1,5 \times 3909 = 20\,767,5$.

На диаграмме размаха для переменной `price` видны два выброса сверху и снизу. На диаграмме размаха переменной `mileage` в нижней части выбросов нет, поэтому нижний «ус» простирается до минимального значения 4867. Вверху мы видим несколько выбросов, превышающих отметку 100 000 миль. Именно эти выбросы ответственны за то, что, как было видно, среднее арифметическое оказалось намного больше, чем медиана.

Визуализация числовых переменных: гистограммы

Гистограмма — это еще одно средство визуализации разброса значений числовой переменной. Гистограмма похожа на диаграмму размаха тем, что делит значения переменной на заранее определенное количество частей, или *интервалов* (bin), которые играют роль «контейнеров» для значений. Однако на этом сходство заканчивается. Если диаграмма размаха состоит из четырех частей, содержащих одинаковое количество значений, но различающихся по диапазону, то на гистограмме используется большее количество частей одинакового диапазона, а ее интервалы могут состоять из разного количества значений.

Для того чтобы построить гистограмму для данных о цене и пробеге подержанных автомобилей, можно воспользоваться функцией `hist()`. Как и в случае с диаграммой размаха, зададим заголовок рисунка, используя параметр `main`, и введем имя оси X с помощью параметра `xlab`. Команды для построения гистограмм выглядят так:

```
> hist(usedcars$price, main = "Histogram of  
Used Car Prices", xlab = "Price ($)")>  
hist(usedcars$mileage, main = "Histogram of Used  
Car Mileage", xlab = "Odometer (mi.)")
```

Гистограммы, которые будут созданы при этом, приведены на рис. 2.2.

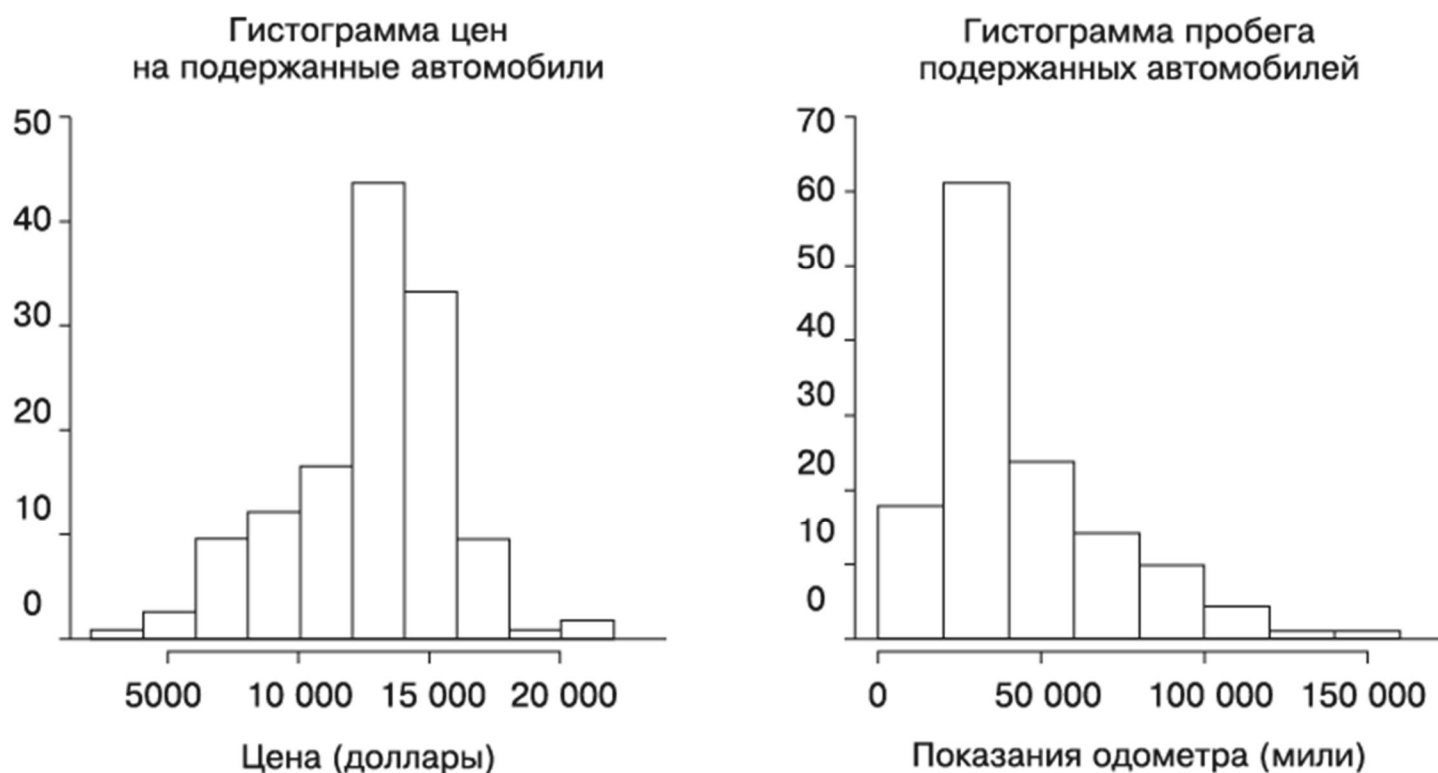


Рис. 2.2. Гистограммы данных о цене и пробеге подержанных автомобилей

Гистограмма состоит из множества столбцов, высота которых соответствует количеству, или *частоте*, значений, попадающих в каждый интервал; все интервалы имеют равную ширину. Вертикальные линии, которые разделяют столбцы, как показано на горизонтальной оси, указывают на начальную и конечную точки диапазона значений, попадающих в данный интервал.



Возможно, вы заметили, что у приведенных выше гистограмм разное количество интервалов. Это связано с тем, что функция `hist()` пытается определить оптимальное количество интервалов для заданного диапазона переменной. Для того чтобы изменить значение, предлагаемое по умолчанию, используйте параметр `breaks`. Если присвоить этому параметру целочисленное значение, например `breaks = 10`, то будет создано ровно десять интервалов одинаковой ширины, а если задать вектор, такой как `c(5000, 10000, 15000, 20000)`, то разделение на интервалы будет выполнено по указанным значениям.

На гистограмме переменной `price` каждый из десяти столбцов соответствует интервалу 2000 долларов, от 2000 до 22 000. Самый высокий столбец, расположенный в центре рисунка, соответствует диапазону от 12 000 до 14 000 долларов и частоте, равной 50. Поскольку известно, что данные охватывают 150 автомобилей, то треть всех автомобилей оценивается на сумму от 12 000 до 14 000 долларов. Почти 90 автомобилей — более половины — продаются по цене от 12 000 до 16 000 долларов.

Гистограмма переменной `mileage` состоит из восьми столбцов, представляющих интервалы по 20 000 миль каждый, от 0 до 160 000 миль. В отличие от гистограммы переменной `price`, самый высокий столбец находится не в центре, а в левой части диаграммы; 70 автомобилей,

попадающих в этот интервал, имеют показания одометра от 20 000 до 40 000 миль.

Вы также могли заметить, что форма этих двух гистограмм несколько различна. Похоже, что цены на подержанные автомобили имеют тенденцию равномерно разделяться по обе стороны от середины, а пробеги автомобилей вытягиваются вправо.

Эта характеристика известна как *коэффициент асимметрии* (skew) или, точнее, *положительный коэффициент асимметрии*, потому что верхние значения (правая сторона) разбросаны гораздо больше, чем нижние (левая сторона). Как показано на рис. 2.3, такие гистограммы данных выглядят растянутыми с одной стороны.

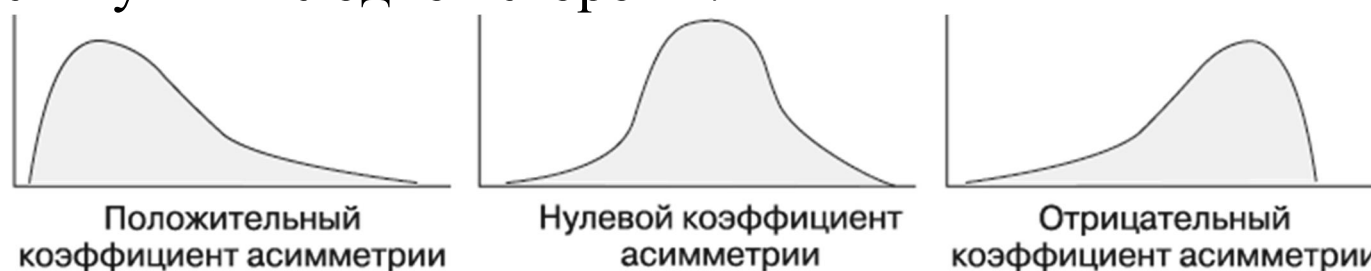


Рис. 2.3. Три примера коэффициента асимметрии, показанные на идеализированных гистограммах

Возможность быстрой диагностики таких паттернов для наших данных является одной из сильных сторон гистограммы как инструмента исследования данных. Это станет еще важнее при изучении других моделей разброса числовых данных.

Интерпретация числовых данных: равномерное и нормальное распределение

Гистограммы, диаграммы размаха и статистические данные, описывающие средние значения и разброс, представляют собой способы исследования распределения значений переменной. *Распределение* переменной описывает вероятность попадания ее значения в различные диапазоны.

Если все значения равновероятны — например, в наборе данных, в котором записаны значения, выпадающие при бросании игрального кубика, — то такое распределение называется *равномерным*. Равномерное распределение легко распознать на гистограмме, поскольку его столбцы примерно одинаковой высоты. Гистограмма равномерного распределения приведена на рис. 2.4.

Важно отметить, что не все случайные события являются равномерно распределенными. Например, бросание шестигранного игрального кубика со смещенным центром тяжести приведет к тому, что одни числа будут выпадать чаще, чем другие. Каждый бросок кубика приводит к случайно выбранному числу, однако вероятность их выпадения неодинакова.



Рис. 2.4. Равномерное распределение, представленное в виде идеализированной гистограммы

Возьмем, к примеру, данные о цене и пробеге подержанных автомобилей. Они явно неоднородны, поскольку одни значения, очевидно, более вероятны, чем другие. В сущности, на гистограмме переменной `price` видно, что значения становятся тем менее вероятными, чем дальше они находятся от обеих сторон центрального столбца, что приводит к колоколообразному распределению. Такое распределение настолько широко распространено среди реальных данных, что является отличительной чертой так называемого *нормального распределения*. Характерная колоколообразная кривая нормального распределения показана на рис. 2.5.



Рис. 2.5. Нормальное распределение, показанное на идеализированной гистограмме

Существует множество типов распределений, не являющихся нормальными, однако многие явления реального мира создают данные, которые описываются с помощью нормального распределения, поэтому свойства нормального распределения изучены очень подробно.

Измерение разброса: дисперсия и стандартное отклонение

Распределение позволяет характеризовать большое количество значений, используя меньшее количество параметров. Нормальное распределение, которое описывает множество типов реальных данных, может быть определено всего двумя параметрами: центром и разбросом. Центр нормального распределения определяется его средним арифметическим значением, которое было использовано ранее. Разброс измеряется статистическим параметром — *стандартным отклонением*.

Для того чтобы рассчитать стандартное отклонение, необходимо сначала вычислить *дисперсию* (variance, Var), которая определяется как среднеквадратическое отклонение между каждым значением и средним арифметическим. В математической нотации дисперсия набора из n значений в наборе с именем X определяется по формуле:

$$\text{Var}(X) = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2.$$

Греческая буква «*мю*» (внешне похожая на *m* или *u*) обозначает среднее арифметическое, а сама дисперсия обозначается греческой буквой «*сигма*» в квадрате.

Стандартное отклонение (standard deviation, StdDev) представляет собой квадратный корень из дисперсии и обозначается буквой «*сигма*»:

$$\text{StdDev}(X) = \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}.$$

В R для вычисления дисперсии и стандартного отклонения используются функции `var()` и `sd()`. Например, если вычислить дисперсию и стандартное отклонение по векторам `price` и `mileage`, то получим следующее:

```
> var(usedcars$price)[1] 9749892>
sd(usedcars$price)[1] 3122.482>
var(usedcars$mileage)[1] 728033954>
sd(usedcars$mileage)[1] 26982.1
```

При интерпретации дисперсии большие числа указывают на то, что данные разбросаны более широко относительно среднего значения. Стандартное отклонение показывает, насколько в среднем каждое значение отличается от среднего арифметического.



Если вычислить эти статистические данные вручную, используя формулы, приведенные выше, то результат будет несколько иной, чем если задействовать стандартные R-функции. Это связано с тем, что в приведенных формулах используется дисперсия совокупности (которая делится на n), а в R — выборочная дисперсия (которая делится на $n - 1$). За исключением очень маленьких наборов данных, различие невелико.

Стандартное отклонение может использоваться для того, чтобы быстро оценить, насколько экстремальным является данное значение, исходя из предположения, что распределение является нормальным. *Правило 68–95–99,7* гласит, что при нормальном распределении 68 % значений находятся в пределах одного стандартного отклонения от среднего арифметического, в то время как 95 и 99,7 % значений находятся в пределах двух и трех стандартных отклонений соответственно. Это показано на рис. 2.6.



Рис. 2.6. Процент значений в пределах одного, двух и трех стандартных отклонений от среднего арифметического при нормальном распределении

Применим эту информацию к данным о подержанных автомобилях. Известно, что среднее арифметическое и стандартное отклонение цены составляли 12 962 и 3122 доллара соответственно. Таким образом, из предположения о нормальном распределении цен следует, что примерно 68 % автомобилей в приведенных данных рекламировались по ценам между $12\,962 - 3122 = 9840$ долларами и $12\,962 + 3122 = 16\,084$ долларами.



Правило 68–95–99,7 применимо только в случае нормального распределения, однако базовый принцип работает для любых данных; значения, превышающие три стандартных отклонения от среднего арифметического, чрезвычайно редки.

Категориальные переменные

Как вы помните, набор данных о подержанных автомобилях содержит три категориальные переменные: `model`, `color` и `transmission`. В R эти переменные сохранены как символьные (`chr`) векторы, а не как тип `factor`, потому что при загрузке данных был использован параметр `stringsAsFactors=FALSE`. Кроме того, хотя переменная `year` хранится в виде числового (`int`) вектора, каждый год можно представить как категорию, применяемую к нескольким автомобилям. Поэтому переменную `year` также можно рассматривать как категориальную.

По контрасту с числовыми данными категориальные данные обычно исследуются с использованием таблиц, а не сводной статистики. Таблица, представляющая одну категориальную переменную, называется *таблицей частотности*. Чтобы построить таблицу частотности для данных о подержанном автомобиле, можно воспользоваться функцией `table()`:

```
> table(usedcars$year)
2000 2001 2002 2003 2004
2005 2006 2007 2008 2009 2010 2011
2012    3    1    1    1    3    2    6   11   14
   42   49   16    1>
table(usedcars$model)
SE    SEL    SES78    23    49>
table(usedcars$color)
Black    Blue    Gold    Gray
Green    Red    Silver    White    Yellow    35    1
7        1    16        5    25        32    16
   3
```

Результатом работы `table()` являются списки категорий номинальной переменной и количество значений, попадающих в каждую категорию. Известно, что набор данных включает сведения о 150 подержанных автомобилях, и можно определить, что примерно треть всех автомобилей выпущена в 2010 году, учитывая, что $49 / 150 = 0,327$.

R также позволяет вычислить пропорции таблицы напрямую, с помощью команды `prop.table()` для таблицы, созданной функцией `table()`:

```
> model_table <- table(usedcars$model)>
prop.table(model_table)          SE          SEL
SES0.5200000      0.1533333      0.3266667
```

Результаты `prop.table()` можно объединять с другими функциями R для представления вывода в другой форме. Предположим, что мы хотим отобразить результаты в процентах с одним десятичным знаком. Для этого можно умножить пропорции на 100, а затем воспользоваться функцией `round()`, указав `digits=1`, как показано в следующем примере:

```
> color_table <- table(usedcars$color)>
color_pct <- prop.table(color_table) * 100>
round(color_pct, digits =
1)Black   Blue   Gold   Gray   Green   Red   Silv
er   White   Yellow23.3   11.3   0.7   10.7
3.3  16.7   21.3   10.7   2.0
```

Это та же информация, что и в результате выполнения `prop.table()` по умолчанию, однако внесенные изменения облегчают чтение данных. Результаты показывают, что наиболее распространены черные машины, в этот цвет окрашена почти четверть (23,3 %) всех рекламируемых автомобилей. На втором месте находится серебристый с показателем 21,3 %, а на третьем — красный, 16,7 %.

Измерение средних значений: мода

В статистической терминологии *мода*— это свойство, значение которого встречается чаще всего. Подобно среднему арифметическому и медиане, мода выступает еще одним показателем среднего значения. Обычно мода используется для категориальных данных, поскольку среднее арифметическое и медиана для номинальных переменных не определены.

Например, в данных о подержанных автомобилях мода переменной `year` равна `2010`, для переменных `model` и `color` — `SE` и `Black` соответственно. Переменная может иметь более одной моды; переменная с одной модой является *унимодальной*, а переменная с двумя

модами — *бимодальной*. Данные с несколькими модами называют *мультимодальными*.



Вы могли подумать, что для вычисления моды в R можно использовать функцию `mode()`, но в R эта функция используется не для вычисления статистической моды, а для получения типа переменной (числовой, списочный и т.п.). Чтобы получить статистическую моду, посмотрите на результат функции `table()` для категории с наибольшим количеством значений.

Одна или несколько мод используются для качественного анализа, позволяя понять важные значения. Однако опасно уделять много внимания модам, поскольку самое распространенное значение не всегда является доминирующим. Например, несмотря на то что черный был единственным самым популярным цветом автомобиля, такой цвет имеют не более четверти всех рекламируемых автомобилей.

Лучше думать о соотношении мод с остальными категориями. Существует ли одна категория, которая доминирует над всеми остальными, или же таких категорий несколько? Подобные размышления о модах способны помочь генерировать проверяемые гипотезы, задавая вопросы о том, почему одни значения более распространены, чем другие. Если черный и серебристый — наиболее популярные цвета подержанных автомобилей, то можно предположить, что здесь описаны автомобили класса люкс, которые, как правило, выпускаются в более консервативных цветах. Или же эти цвета могут указывать на автомобили экономкласса, выбор цвета которых весьма ограничен. Запомним это и продолжим рассматривать данные.

Считая моды стандартными значениями, мы можем применить концепцию статистической моды к числовым данным. Строго говоря, едва ли мода будет иметь смысл для непрерывной переменной, так как у нее не существует двух повторяющихся значений. Однако, если представить моду как самый высокий столбец на гистограмме, можно исследовать моды таких переменных, как `price` и `mileage`. Полезно будет изучить моду при исследовании числовых данных, в частности, чтобы проверить, являются ли данные мультимодальными (рис. 2.7).

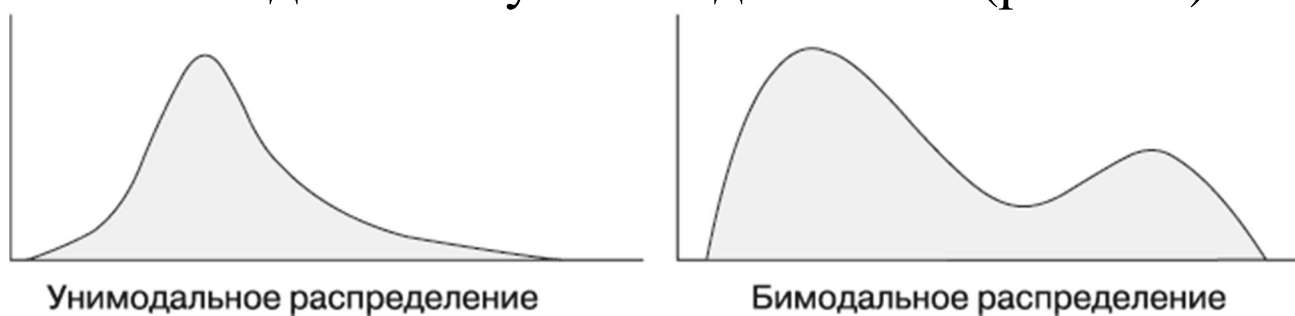


Рис. 2.7. Гипотетические распределения числовых данных с одной и двумя модами

Взаимосвязи между переменными

До сих пор мы рассматривали переменные по отдельности, вычисляя только *одномерную статистику*. В процессе были заданы вопросы, на которые вы не могли ответить раньше.

- Означают ли данные о цене и пробеге, что исследуются только автомобили экономкласса, или среди данных есть автомобили класса люкс с большим пробегом?
- Можно ли на основании взаимосвязей между моделью и цветом сделать вывод о типе автомобилей, которые мы исследуем?

Ответы на подобные вопросы можно получить путем изучения *двумерных отношений*, то есть взаимосвязей между двумя переменными. Отношения между более чем двумя переменными называются *многомерными*. Начнем с двумерного случая.

Визуализация отношений: диаграммы разброса

Диаграмма разброса визуализирует двумерные отношения между числовыми признаками. Это двумерное изображение, где на координатной плоскости нанесены точки, так что значения одного признака откладываются по горизонтальной оси X , а значения другого признака — по вертикальной оси Y . Паттерны расположения точек отражают основные взаимосвязи между этими двумя признаками.

Чтобы ответить на вопрос о соотношении цены и пробега, построим диаграмму разброса. Для этого воспользуемся функцией `plot()` с параметрами `main`, `xlab` и `ylab`.

Чтобы применить `plot()`, нам нужно определить векторы `x` и `y`, содержащие значения, используемые для позиционирования точек на рисунке. Хотя выводы не зависят от того, какая именно переменная будет выбрана для оси X , а какая — для оси Y , общее соглашение таково, что в качестве переменной `y` выбирается та, которая предположительно зависит от другой (и поэтому называется *зависимой переменной*). Поскольку продавец не может изменить показания одометра автомобиля, пробег вряд ли будет зависеть от цены. Наоборот, гипотеза состоит в том, что цена автомобиля зависит от показаний одометра. Поэтому в качестве зависимой переменной `y` мы выберем цену.

Полностью команда построения диаграммы разброса выглядит так:

```
> plot(x = usedcars$mileage, y =  
usedcars$price, main = "Scatterplot of  
Price vs. Mileage", xlab = "Used Car  
Odometer (mi.)", ylab = "Used Car Price  
($)")
```

В результате получим диаграмму разброса, представленную на рис. 2.8.

На данной диаграмме разброса четко видна зависимость между ценой на подержанный автомобиль и показаниями одометра. Чтобы прочесть диаграмму, изучите, как изменяются значения переменной **Y** при увеличении значений на оси **X**. В этом случае, как правило, чем больше пробег, тем ниже цена автомобиля. Если вам когда-либо приходилось продавать или покупать подержанный автомобиль, то вряд ли это будет для вас открытием.

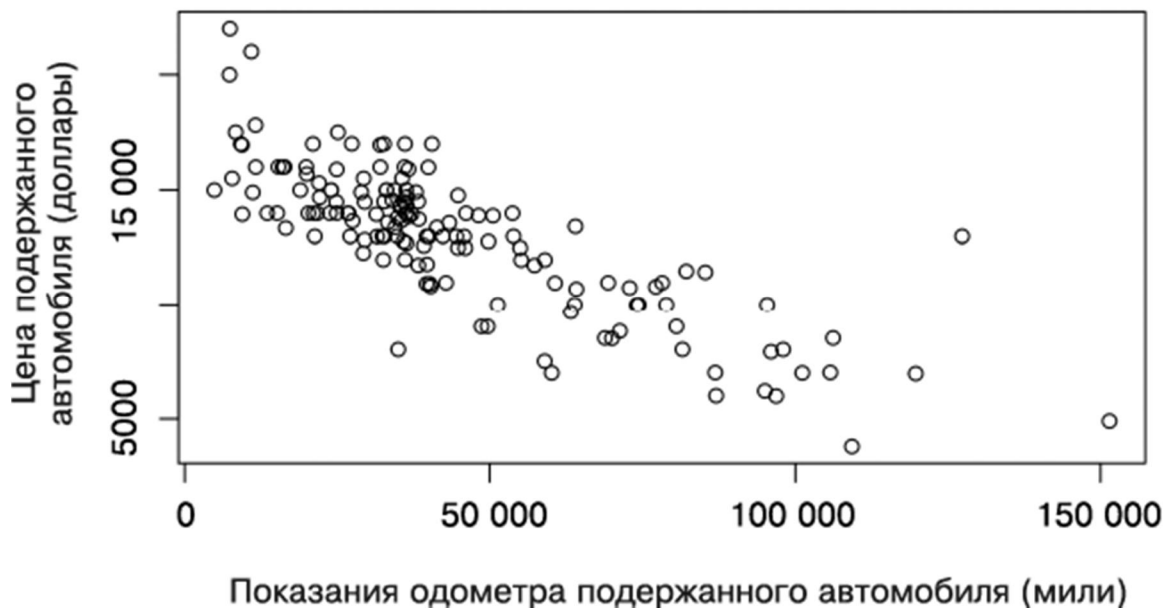


Рис. 2.8. Зависимость цены подержанного автомобиля от пробега

Возможно, более интересен тот факт, что очень мало автомобилей имеют и высокую цену, и большой пробег, за исключением одинокого выброса на отметке 125 000 миль и 14 000 долларов. Отсутствие большего количества точек, подобных этой, свидетельствует о том, что этот набор данных вряд ли включает в себя автомобили класса люкс с большим пробегом.

Все самые дорогие автомобили в наборе данных, особенно те, что стоят свыше 17 500 долларов, имеют очень низкий пробег. Это говорит о том, что мы имеем дело с автомобилями одного типа, которые, будучи новыми, продаются примерно за 20 000 долларов.

Отношение, которое наблюдается, между ценами на автомобили и пробегом, — отрицательная связь, она образует паттерн точек, представляющий собой линию, направленную вниз. Аналогично положительная связь образует линию, направленную вверх. Ровная линия, или кажущийся случайным разброс точек, свидетельствует о том, что эти две переменные вообще не связаны. Сила линейной связи между двумя переменными измеряется величиной, известной как *корреляция*. Подробнее о корреляции читайте в главе 6, где рассматриваются методы моделирования линейных отношений.



Имейте в виду, что не все связи образуют прямые линии. Иногда точки образуют U- или V-образную форму, а в других случаях паттерн значения

переменных x или y кажется более или менее ярко выраженным. Такие модели подразумевают, что взаимосвязи между переменными не являются линейными.

Исследование взаимосвязей: перекрестные таблицы

Для изучения взаимосвязи между двумя номинальными переменными используется *перекрестная таблица* (также известная как *кросс-таблица* или *таблица сопряженности*). Кросс-таблица похожа на диаграмму разброса тем, что позволяет исследовать, каким образом значения одной переменной изменяются в зависимости от значений другой. В этой таблице строки представляют собой уровни одной переменной, а столбцы — уровни другой переменной. Число в каждой ячейке таблицы показывает количество значений, попадающих в конкретное пересечение строки и столбца.

Чтобы ответить на вопрос о том, существует ли связь между моделью и цветом, рассмотрим кросс-таблицу. В R есть несколько функций для создания перекрестных таблиц, включая `table()`, которую мы использовали для формирования таблиц частотности. Пожалуй, самой удобной является функция `CrossTable()` из пакета `gmodels`, созданного Грегори Р. Уорнсом (Gregory R. Warnes), так как она отображает в одной таблице проценты строк, столбцов и полей, избавляя нас от необходимости вычислять их самостоятельно. Чтобы установить пакет `gmodels`, введите следующую команду:

```
> install.packages("gmodels")
```

После установки пакета введите `library(gmodels)`, чтобы загрузить пакет. Это необходимо делать в каждом сеансе R, в котором планируется использовать функцию `CrossTable()`.

Прежде чем приступить к анализу, упростим наш проект, сократив количество уровней в переменной `color`. Эта переменная имеет девять уровней, но нам не нужно так много подробностей. На самом деле нас интересует, относится ли цвет машины к консервативным. С этой целью мы разделим девять цветов на две группы. Первая группа будет включать в себя консервативные цвета: `Black`, `Gray`, `Silver` и `White`; вторая — `Blue`, `Gold`, `Green`, `Red` и `Yellow`. Создадим двоичную индикаторную переменную (часто называемую *фиктивной переменной*), указывающую, относится ли цвет машины к консервативным, согласно нашему определению. Эта переменная будет принимать значение `1`, если цвет консервативный, и `0` — если нет:

```
> usedcars$conservative <- usedcars$color %in%  
c("Black", "Gray", "Silver", "White")
```

Возможно, вы заметили новую команду:

оператор `%in%` возвращает `TRUE` или `FALSE` для каждого значения вектора, стоящего с левой стороны от оператора, указывая, найдено ли в векторе значение, заданное с правой стороны. Проще говоря, эту строку можно истолковать так: «Подержанный автомобиль черного, серого, серебристого или белого цвета?»

Исследуя результат выполнения функции `table()` для нашей новой переменной, можно увидеть, что примерно две трети автомобилей имеют консервативные цвета, а одна треть — нет:

```
>
```

```
table(usedcars$conservative) FALSE TRUE 51
99
```

Теперь посмотрим на кросс-таблицу, чтобы увидеть, как количество автомобилей консервативных цветов зависит от модели. Поскольку, предположительно, модель автомобиля диктует выбор цвета, рассмотрим индикатор консервативного цвета как зависимую переменную (y).

Поэтому команда `CrossTable()` будет выглядеть так:

```
> CrossTable(x = usedcars$model, y =
usedcars$conservative)
```

Результат будет следующим:

Cell Contents

	usedcars\$conservative		
usedcars\$model	FALSE	TRUE	Row Total
SE	27	51	78
	0.009	0.004	
	0.346	0.654	0.520
	0.529	0.515	
	0.180	0.340	
SEL	7	16	23
	0.086	0.044	
	0.304	0.696	0.153
	0.137	0.162	
	0.047	0.107	
SES	17	32	49
	0.007	0.004	
	0.347	0.653	0.327
	0.333	0.323	
	0.113	0.213	
Column Total	51	99	150
	0.340	0.660	

В результатах функции `CrossTable()` содержится множество данных. Расположенная сверху легенда (с заголовком `CellContents`) указывает на то, как следует интерпретировать каждое значение. В строках таблицы указаны три модели подержанных автомобилей: `SE`, `SEL` и `SES` (и еще одна дополнительная строка для всех моделей). Столбцы указывают, относится ли цвет машины к консервативным (плюс еще один столбец для обоих типов цветов).

Первое значение в каждой ячейке указывает, какое количество автомобилей имеет данную комбинацию модели и цвета. Пропорции показывают вклад каждой ячейки в статистику хи-квадрат, сумму строк, сумму столбцов и общую сумму таблицы.

Наибольший интерес вызывает доля консервативных автомобилей для каждой модели. Пропорции строк говорят о том, что 0,654 (65 %) автомобилей модели `SE` окрашены в консервативные цвета, по сравнению с 0,696 (70 %) автомобилей модели `SEL` и 0,653 (65 %) модели `SES`. Эти различия относительно невелики, что свидетельствует об отсутствии существенных различий в типах цветов, выбранных для каждой модели автомобиля.

Значения хи-квадрат означают вклад данной ячейки в *критерий согласия Пирсона*, или *критерий согласия хи-квадрат* для двух переменных. Данный критерий показывает, насколько вероятным является то, что разница в количестве ячеек в таблице обусловлена одной лишь случайностью. Если вероятность очень низкая, это дает убедительные доказательства того, что две переменные связаны.

Для того чтобы вычислить критерий хи-квадрат, можно добавить в вызов функции `CrossTable()` дополнительный параметр: `chisq=TRUE`. В данном случае вероятность составляет около 93 %, а это говорит о том, что, весьма вероятно, различия в количестве ячеек обусловлены только случайностью, а не действительной взаимосвязью между моделью и цветом.

Резюме

Из этой главы вы узнали об основах управления данными в R. Начали с глубокого взгляда на структуры, используемые для хранения различных типов данных. Основной структурой данных R является вектор, который расширяется и объединяется в более сложные типы данных, такие как списки и фреймы данных. Фрейм представляет собой структуру данных R, которая соответствует понятию набора данных, имеющего признаки и примеры. В R есть функции для чтения и записи фреймов данных в файлы табличных данных, такие как электронные таблицы.

Затем мы рассмотрели реальный набор данных, содержащий цены на подержанные автомобили. Мы исследовали числовые переменные, используя сводную статистику в виде среднего арифметического и разброса и визуализировали взаимосвязь между ценами и показаниями одометра с помощью диаграммы разброса. Далее с помощью таблиц мы проанализировали номинальные переменные. В ходе изучения данных о подержанных автомобилях мы провели исследование, которое можно использовать для изучения любого набора данных. Эти навыки потребуются и для других проектов, которые будут рассмотрены в книге.

Теперь, усвоив основы управления данными с помощью R, вы готовы начать использовать машинное обучение для решения реальных задач. В следующей главе вам предстоит решить первую задачу классификации, используя метод ближайших соседей.

3. Ленивое обучение: классификация с использованием метода ближайших соседей

Сейчас в разных городах мира набирает популярность необычный способ пообедать в ресторане. В абсолютно темном зале посетителей обслуживают официанты, которые двигаются по заранее изученным маршрутам, ориентируясь только на осязание и слух. Вся прелесть этих заведений заключается в том, что многие убеждены: лишив себя возможности видеть, можно обострить чувства вкуса и обоняния, и пища будет восприниматься по-новому.

Можете ли вы представить, как посетитель воспринимает на вкус еду, которую не видит? Уже с первого кусочка чувства обостряются. Каковы доминирующие ароматы? Какова пища на вкус: острая или сладкая? Похоже ли это на что-то, что вы ели раньше? Я ассоциирую этот исследовательский процесс со слегка видоизмененной пословицей: если что-то пахнет как утка и имеет вкус утки, то вы, вероятно, едите утку.

Такая же идея может быть использована при машинном обучении. Пословица в тему: рыбак рыбака видит издалека. Другими словами, вещи, похожие друг на друга, скорее всего, будут иметь одинаковые свойства. В ML этот принцип используется для классификации данных: данные помещаются в одну категорию с аналогичными или ближайшими соседями. Эта глава посвящена классификаторам, в которых используется такой подход. Мы рассмотрим:

- ключевые понятия, которые определяют методы ближайших соседей, и то, почему эти алгоритмы называются ленивым обучением;
- методы измерения сходства двух объектов с помощью расстояния;
- применение популярного метода ближайших соседей, который называется k-NN.

Если все эти разговоры о еде пробудили ваш аппетит, то нашей первой задачей будет изучить метод k -NN, применив его к давним кулинарным спорам.

Что такое классификация методом ближайших соседей

В двух словах, методы *ближайших соседей* определяются по их свойству классифицировать немаркированные объекты путем присвоения им класса похожих маркированных объектов. По аналогии с примером ужина в темноте, описанным выше, когда человек идентифицирует новую еду, сравнивая ее с тем, что он встречал ранее, при классификации методом ближайших соседей компьютеры демонстрируют почти человеческую способность вспоминать полученный опыт, чтобы делать выводы о текущих обстоятельствах. Несмотря на простоту этой идеи, методы ближайших соседей являются чрезвычайно мощными. Они успешно используются для:

- приложений компьютерного зрения, включая оптическое распознавание символов и распознавание лиц как на неподвижных изображениях, так и на видео;
- систем выдачи рекомендаций, которые прогнозируют, понравится ли данному человеку фильм или песня;
- выявления закономерностей в генетических данных для выделения конкретного белкового продукта или определения заболеваний.

В целом метод ближайших соседей хорошо подходит для задач классификации, где взаимосвязи между признаками и целевыми классами многочисленны, сложны или чрезвычайно трудны для понимания, хотя элементы схожего типа классов имеют тенденцию быть однородными. Другими словами, если концепцию сложно определить, но вы узнаете ее, когда увидите, то метод ближайших соседей может быть уместен. Однако если данные зашумлены и, следовательно, нет четкого различия между группами, то алгоритмы ближайших соседей могут с трудом определить границы классов.

Алгоритм k -NN

Классификация методом ближайших соседей иллюстрируется алгоритмом *ближайших соседей* (k -NN). Это, возможно, один из самых простых алгоритмов машинного обучения, тем не менее он широко используется. Преимущества и недостатки этого алгоритма представлены в табл. 3.1.

Алгоритм k -NN получил свое название из-за того, что в нем для классификации еще не отнесенных к какому-либо классу объектов используется информация о его k ближайших соседях. Буква k является переменной, означающей, что можно использовать любое количество ближайших соседей. После выбора k алгоритму нужен тренировочный

набор данных, состоящий из объектов, которые уже были классифицированы на несколько категорий, помеченных номинальной переменной. Затем для каждого неклассифицированного объекта из тестового набора данных k -NN выбирает k объектов в тренировочном наборе данных, которые по принципу подобия являются к нему «ближайшими». Неклассифицированному тестовому объекту присваивается класс, к которому относится большинство из k ближайших соседей.

Таблица 3.1

Преимущества	Недостатки
Простота и эффективность. Не дает никаких предположений о базовом распределении данных. Быстрая фаза обучения	Не генерирует модель, ограничивая способность понимать, каким образом признаки связаны с классом. Требуется выбор подходящего k . Медленная фаза классификации. Номинальные признаки и отсутствующие данные требуют дополнительной обработки

Чтобы наглядно объяснить процесс, вернемся к опыту слепой дегустации, описанному выше. Предположим, что, перед тем, как съесть неизвестную пищу, мы создали набор данных с записями наших впечатлений от нескольких продуктов, которые пробовали ранее. Для простоты мы оценили только два признака каждого продукта. Первый — это показатель степени хрусткости, от 1 до 10, а второй — степень сладости, тоже от 1 до 10 баллов. Затем мы поделили все продукты на три вида: фрукты, овощи и белковые продукты, игнорируя пищу, насыщенную углеводами и жирами. Первые несколько строк такого набора данных могут быть структурированы следующим образом (табл. 3.2).

Таблица 3.2

Продукт	Сладость	Хрусткость	Тип продукта
Яблоко	10	9	Фрукт
Бекон	1	4	Белковый продукт
Банан	10	1	Фрукт
Морковь	7	10	Овощ
Сельдерей	3	10	Овощ
Сыр	1	1	Белковый продукт

Алгоритм k -NN рассматривает объекты как координаты в многомерном пространстве признаков. Поскольку набор данных о продуктах включает в себя только два признака, то его пространство признаков является двумерным. Мы можем нанести двумерные данные на диаграмму разброса, где значение x обозначает сладость продукта, а значение y — хрусткость. После добавления еще нескольких продуктов в набор данных вкусов диаграмма разброса может выглядеть, как на рис. 3.1.

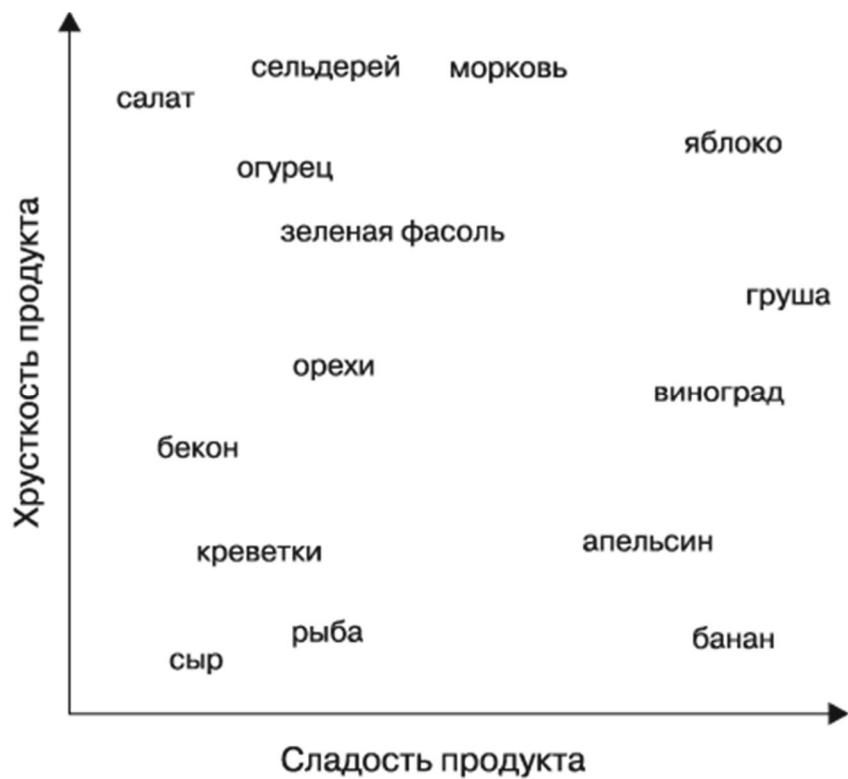


Рис. 3.1. Диаграмма разброса хрустящих и сладких продуктов

Заметили тенденцию? Продукты одного вида, как правило, расположены близко друг от друга. Как видно на рис. 3.2, овощи имеют тенденцию быть хрустящими, но не сладкими; фрукты имеют тенденцию быть сладкими и либо хрустящими, либо нет; белковые продукты, как правило, не являются ни хрустящими, ни сладкими.

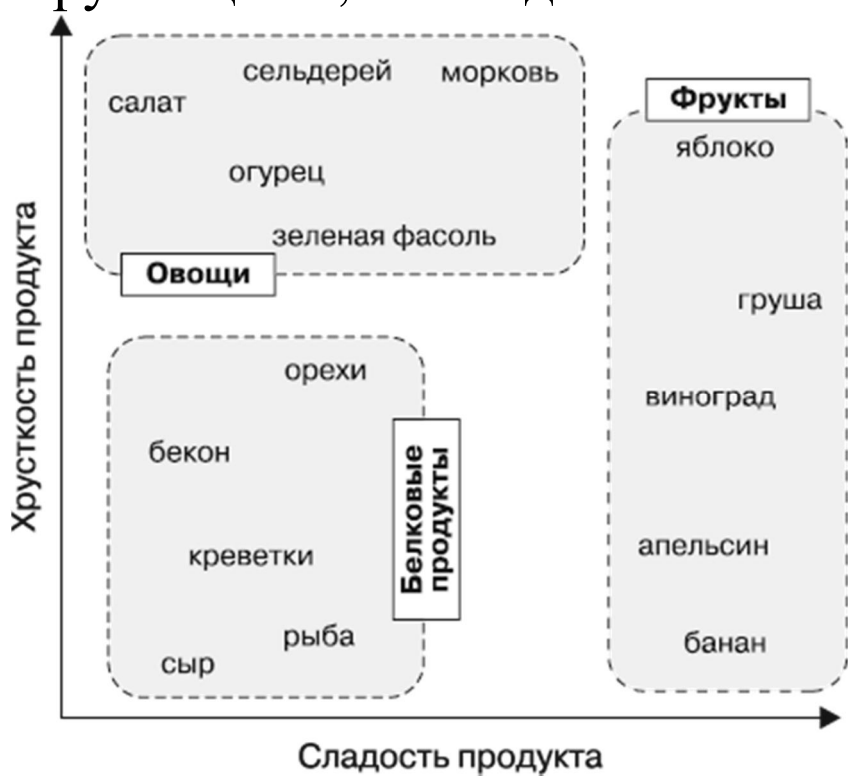


Рис. 3.2. Продукты одинакового типа имеют сходные свойства

Предположим, что после формирования этого набора данных мы решили использовать его для решения вечного вопроса: помидор — это фрукт или овощ? Мы можем использовать метод ближайших соседей, чтобы определить, какой класс лучше подходит, как показано на рис. 3.3.

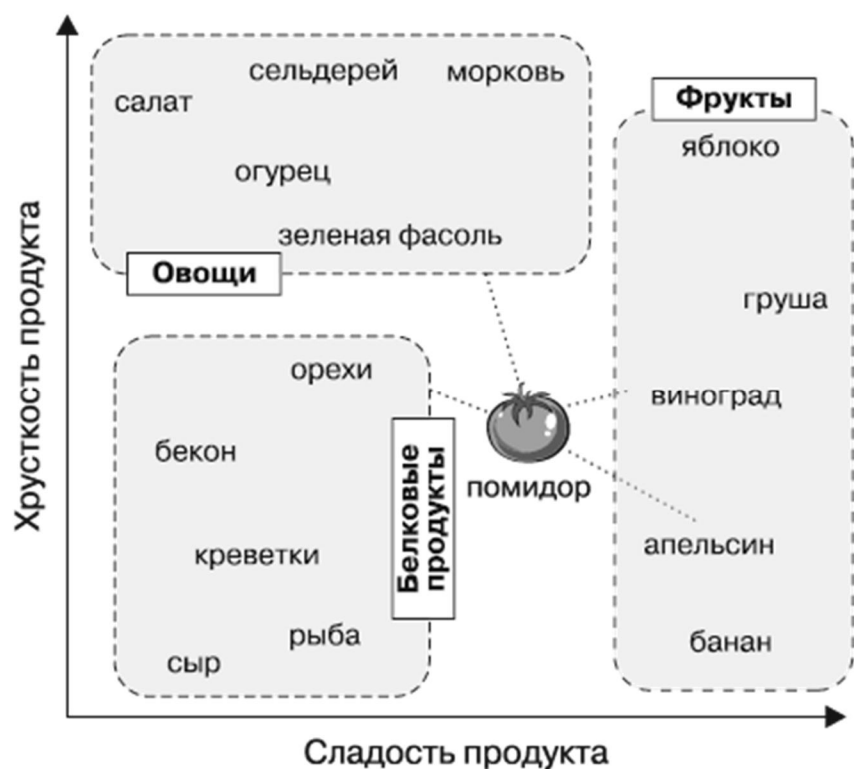


Рис. 3.3. Ближайшие соседи помидора дают представление о том, фрукт это или овощ

Измерение степени сходства с помощью расстояния

Для нахождения ближайших соседей помидора нужна *функция расстояния*. Она представляет собой формулу, по которой можно измерить степень сходства между двумя экземплярами.

Существует множество способов вычисления расстояния. Традиционно алгоритм k-NN использует *евклидово расстояние* — расстояние, которое можно измерить, если соединить две точки линейкой. Эти расстояния показаны на рис. 3.3 пунктирными линиями, соединяющими помидор с соседними объектами.



Евклидово расстояние измеряется «по прямой линии», подразумевая кратчайший прямой маршрут. Другим распространенным показателем расстояния является *манхэттенское расстояние* — путь, по которому пешеход идет по городским кварталам. Подробнее о других способах измерения расстояния читайте в документации по функции расстояния R с помощью команды `?dist`.

Евклидово расстояние определяется по следующей формуле, где p и q — объекты, которые нужно сравнить и каждый из которых имеет n признаков. Термин p_1 обозначает значение первого признака объекта p , q_1 — значение первого признака объекта q :

$$\text{dist}(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}.$$

Формула расстояния подразумевает сравнение значений признаков каждого объекта. Например, для того, чтобы вычислить расстояние между помидором (сладость = 6, хрусткость = 4) и спаржей (сладость = 3, хрусткость = 7), можно использовать следующую формулу:

$$\text{dist}(\text{помидор}, \text{спаржа}) = \sqrt{(6 - 3)^2 + (4 - 7)^2} = 4,2.$$

Аналогичным образом можно вычислить расстояние между помидором и еще несколькими ближайшими соседями (табл. 3.3).

Таблица 3.3

Продукт	Сладость	Хрусткость	Тип еды	Расстояние до помидора
Виноград	8	5	Фрукт	$\sqrt{(6 - 8)^2 + (4 - 5)^2} = 2,2$
Спаржа	3	7	Овощ	$\sqrt{(6 - 3)^2 + (4 - 7)^2} = 4,2$
Орехи	3	6	Белковый продукт	$\sqrt{(6 - 3)^2 + (4 - 6)^2} = 3,6$
Апельсин	7	3	Фрукт	$\sqrt{(6 - 7)^2 + (4 - 3)^2} = 1,4$

Чтобы классифицировать помидор как овощ, фрукт или белковый продукт, сначала назначим ему тип продукта его ближайшего соседа. Это называется классификацией 1-NN, потому что $k = 1$. Единственный ближайший сосед помидора с расстоянием 1,4 — это апельсин. Поскольку апельсин — фрукт, то алгоритм 1-NN классифицирует помидор как фрукт.

Если же использовать алгоритм k -NN с $k = 3$, то будет проведено голосование среди трех ближайших соседей: апельсина, винограда и ореха. Поскольку большинство из этих соседей — фрукты (два из трех голосов), то помидор снова классифицируется как фрукт.

Выбор подходящего k

Решение о том, сколько соседей использовать для алгоритма k -NN, определяет, насколько хорошо модель будет обобщена для будущих данных. Проблема поиска баланса между перетренированностью и недотренированностью модели называется *компромиссом между смещением и разбросом*. При слишком большом k уменьшается разброс, или дисперсия (variance), вызванная зашумленными данными, но это может привести к смещению (bias) так, что обучаемый рискует игнорировать небольшие, но важные паттерны.

Предположим, что мы заняли крайнюю позицию, установив очень большое значение k , равное общему количеству наблюдений в тренировочных данных. Тогда при появлении каждого нового обучающего экземпляра, представленного в итоговом голосовании, самый распространенный класс всегда будет иметь большинство избирателей. Следовательно, модель постоянно будет прогнозировать класс большинства, независимо от ближайших соседей.

С другой стороны, использование единственного ближайшего соседа приводит к тому, что зашумление данных и выбросы могут сильно влиять на классификацию объектов. Например, предположим, что некоторые из обучающих объектов были маркированы случайно. Тогда любой немаркированный объект, ближайший к неправильно маркированному соседу, будет иметь неправильный класс, даже если девять других соседей проголосовали бы по-другому.

Очевидно, что наилучшее значение k находится где-то посередине между этими крайностями.

На рис. 3.4 в обобщенном виде показано, каким образом большие и меньшие значения k влияют на границу решения (изображена пунктирной линией). Меньшие значения допускают более сложные границы принятия решений, которые более точно соответствуют тренировочным данным. Проблема в том, что мы не знаем, какая граница — прямая или изогнутая — находится ближе к истинной концепции, которую следует изучить.

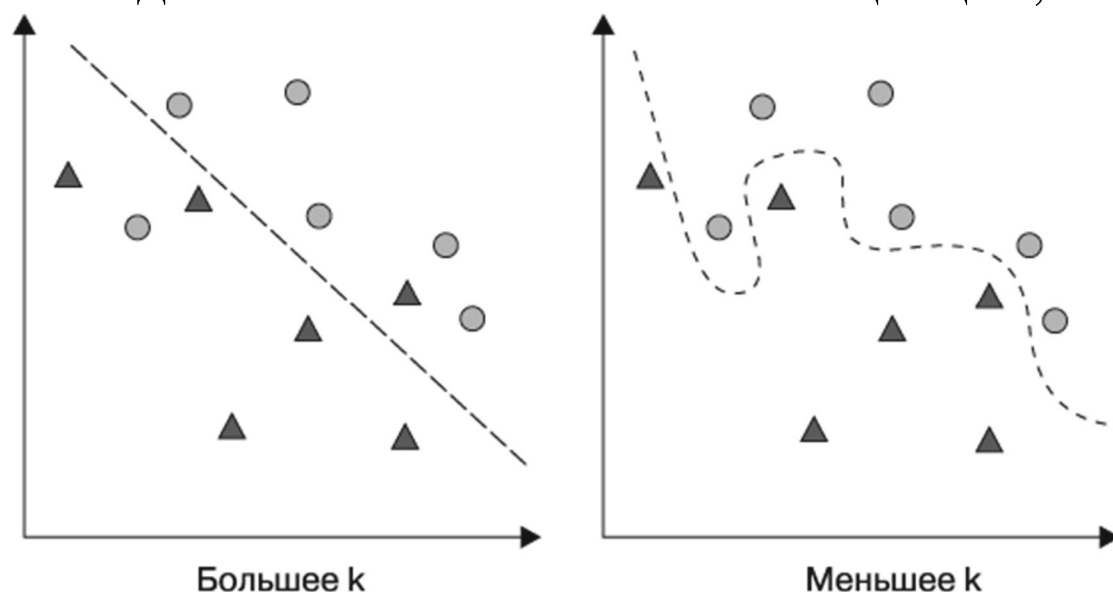


Рис. 3.4. Большее значение k имеет более высокое смещение и меньшую дисперсию, чем меньшее k

На практике выбор k зависит от сложности изучаемой концепции и количества записей в тренировочных данных. Одна из распространенных методик — начать с k , равного квадратному корню из числа обучающих объектов. В классификаторе продуктов питания, который мы построили, можно установить $k = 4$, потому что в тренировочных данных было 15 примеров продуктов, а квадратный корень из 15 равен 3,87.

Однако такие правила не всегда приводят к единственному наилучшему k . Другой подход состоит в том, чтобы протестировать несколько значений k на различных тестовых наборах данных и выбрать тот, который обеспечивает наилучшую эффективность классификации. Тем не менее если данные не очень зашумленные, то большой тренировочный набор данных может снизить важность выбора k . Это связано с тем, что даже у точных наборов данных будет достаточно большое количество объектов, чтобы учесть голоса ближайших соседей.



Менее распространенное, но интересное решение этой проблемы — выбрать большее k и использовать взвешенное голосование, при котором голоса ближайших соседей считаются более весомыми, чем голоса соседей, расположенных далеко. Некоторые реализации k -NN предполагают такой вариант.

Подготовка данных для использования в алгоритме k -NN

Перед применением алгоритма k-NN признаки обычно подвергаются нормировке. Основанием для этого является то, что формула расстояния сильно зависит от способа измерения объектов. В частности, если одни объекты имеют гораздо больший диапазон значений, чем другие, то при вычислении расстояний будут преобладать элементы с более широкими диапазонами. В примере с дегустацией пищи это не было проблемой, так как сладость и хрусткость измерялись по шкале от 1 до 10.

Однако предположим, что мы добавили в набор данных еще один признак, описывающий остроту пищи, которая измеряется по шкале Сковилла. Для тех, кто не знаком с этим показателем, поясню: это стандартная мера остроты специй, варьирующаяся от нуля (совсем не острая) до миллиона и более (для самых острых сортов перца чили). Поскольку разница между острыми и неострыми продуктами может превышать миллион, а разница между сладкой и несладкой или хрустящей и нехрустящей едой составляет не более 10, то разница в масштабе позволит уровню остроты влиять на функцию расстояния намного сильнее, чем два других фактора. Не корректируя данные, мы можем увидеть, что показатели расстояния различают продукты только по их остроте; влияние хрусткости и сладости будет уменьшено.

Решение состоит в том, чтобы изменить масштабы признаков, сузив или расширив их диапазон так, чтобы каждый вносил в формулу расстояния относительно равный вклад. Например, если сладость и хрусткость измеряются по шкале от 1 до 10, то можно было бы измерять остроту также по шкале от 1 до 10. Есть несколько распространенных способов выполнить такое масштабирование.

Традиционным методом изменения масштаба для алгоритма k-NN является *минимаксная нормализация*. В процессе такой нормализации признак преобразуется так, что все его значения попадают в диапазон от 0 до 1. Формула нормализации признака выглядит так:

$$X_{\text{new}} = \frac{X - \min(X)}{\max(X) - \min(X)}.$$

В сущности, в этой формуле из каждого значения признака вычитается минимальное значение x , и результат делится на диапазон x .

Полученные нормализованные значения признаков можно интерпретировать как указание на то, насколько далеко, в диапазоне от 0 до 100 %, исходное значение отстоит от исходного минимума и максимума.

Еще одно распространенное преобразование называется *стандартизацией по z-оценке*. В следующей формуле из значения признака x вычитается среднее арифметическое и результат делится на стандартное отклонение x :

$$X_{\text{new}} = \frac{X - \mu}{\sigma} = \frac{X - \text{Mean}(X)}{\text{StdDev}(X)}.$$

Эта формула, основанная на свойствах нормального распределения, описанного в главе 2, изменяет масштаб каждого значения признака с точки зрения того, на сколько стандартных отклонений эти значения больше или меньше среднего арифметического значения. Полученное значение называется *z-оценкой*. Z-оценки образуют неограниченный диапазон отрицательных и положительных чисел. В отличие от нормализованных значений, у них нет predetermined минимума и максимума.



К тестовым примерам, которые впоследствии будут классифицироваться по алгоритму k-NN, должен применяться тот же метод масштабирования, который был использован для тренировочного набора данных. В случае минимаксной нормализации это может привести к сложной ситуации, так как минимум или максимум для будущих примеров может находиться вне диапазона значений, наблюдаемых в тренировочных данных. Если заранее известно правдоподобное минимальное или максимальное значение, то можно использовать эти константы вместо наблюдаемых значений. В качестве альтернативы можно использовать стандартизацию по z-оценке, исходя из предположения, что будущие примеры будут иметь то же среднее арифметическое и стандартное отклонение, что и тренировочные. Для номинальных данных евклидова формула расстояния не определена. Поэтому, для того чтобы вычислить расстояние между номинальными объектами, нужно преобразовать их в числовой формат. Типичное решение предполагает *фиктивное кодирование*, в котором значение 1 указывает на одну категорию, а 0 — на другую. Например, фиктивное кодирование для гендерной переменной может быть построено так:

$$\text{male} = \begin{cases} 1 & \text{если } x = \text{male}, \\ 0 & \text{в иных случаях.} \end{cases}$$

Обратите внимание, что фиктивное кодирование двухкатегорийной (двоичной) гендерной переменной приводит к появлению единственной новой функции с именем `male`. Нет необходимости создавать отдельную функцию для женщины. Поскольку оба пола являются взаимоисключающими, достаточно знать один из них.

Это верно и при обобщении. Номинальный n -категорийный признак можно фиктивно закодировать, создавая двоичные переменные для $n - 1$ уровней этого признака. Например, фиктивное кодирование для трехкатегорийной температурной переменной — например, принимающей значения «горячо» (`hot`), «средне» (`medium`) или «холодно» (`cold`) — может быть представлено в виде $(3 - 1) = 2$ признаков:

$$\text{hot} = \begin{cases} 1 & \text{если } x = \text{hot}, \\ 0 & \text{в иных случаях;} \end{cases}$$

$$\text{medium} = \begin{cases} 1 & \text{если } x = \text{medium} \\ 0 & \text{в иных случаях.} \end{cases}$$

Зная, что оба показателя, **hot** и **medium**, равны **0**, можно сделать вывод о том, что температура холодная. Нам не нужна третья функция для категории «холодно». Поскольку только один атрибут закодирован как **1**, а все остальные должны быть равны **0**, то фиктивное кодирование также известно как *прямое кодирование*.

Удобной особенностью фиктивного кодирования является то, что расстояние между фиктивными объектами всегда равно 1 или 0, и, таким образом, значения находятся в том же масштабе, что и числовые данные, полученные путем минимаксной нормализации. Никаких дополнительных преобразований не требуется.



Если номинальный признак является упорядоченным (как в случае с температурой), то вместо фиктивного кодирования можно использовать нумерацию категорий и применить нормализацию. Например, состояния «холодно», «средне» и «горячо» можно пронумеровать как 1, 2 и 3, что нормализуется до 0, 0,5 и 1. Однако необходимо учитывать, что данный подход следует применять только в том случае, если расстояния между категориями эквивалентны. Например, категории доходов для бедных, среднего класса и богатых упорядочены, однако разница между бедным и средним классом может отличаться от разницы между средним классом и богатым. Если расстояния между группами не равны, то более безопасным подходом будет фиктивное кодирование.

Почему алгоритм k-NN называют ленивым

Алгоритмы классификации, основанные на методах ближайших соседей, называют *ленивыми алгоритмами обучения*, потому что абстрагирования при этом не происходит. Процессы абстрагирования и обобщения полностью пропускаются, что нарушает определение обучения, предложенное в главе 1.

Согласно строгому определению обучения ленивый обучаемый на самом деле ничему не учится. Он лишь дословно запоминает тренировочные данные. Это позволяет очень быстро проходить этап обучения, на котором фактически ничто не выучивается.

Конечно, недостатком такого метода является то, что прогнозирование выполняется сравнительно медленно. Вследствие сильной зависимости от тренировочных примеров, а не от абстрактной модели ленивое обучение также известно как *экземпляр-ориентированное обучение* (instance-based learning), ИЛИ *обучение методом заучивания наизусть*.

Поскольку обучаемые на основе экземпляров не строят модель, говорят, что этот метод относится к классу *непараметрических* методов обучения — никакие параметры данных не изучаются. Не создавая теории о базовых данных, непараметрические методы ограничивают нашу способность понимать, каким образом классификатор использует данные. С другой стороны, это позволяет обучаемому найти естественные закономерности, вместо того чтобы пытаться вписать данные в предвзятую и потенциально неверную функциональную форму (рис. 3.5).



Рис. 3.5. Разные алгоритмы машинного обучения имеют разные смещения и могут прийти к разным выводам!

Несмотря на то что классификаторы k -NN являются ленивыми, они весьма мощные. Скоро вы увидите, как простые принципы обучения методом ближайших соседей можно использовать для автоматизации процесса скрининга рака.

Пример: диагностика рака молочной железы с помощью алгоритма k -NN

Обычный скрининг рака молочной железы позволяет диагностировать и вылечить это заболевание до того, как появятся заметные симптомы. Процесс раннего выявления включает в себя исследование ткани молочной железы на наличие аномальных уплотнений или новообразований. Если такое уплотнение обнаружится, то выполняется аспирационная биопсия с использованием полой тонкой иглы, которой из этого новообразования извлекают небольшое количество клеток. Затем врач рассматривает клетки под микроскопом и определяет, злокачественное это новообразование или доброкачественное.

Если бы ML могло автоматизировать идентификацию раковых клеток, это принесло бы значительную пользу системе здравоохранения. Автоматизированные процессы, очевидно, повысят эффективность процесса выявления рака, что сократит время диагностики и позволит уделять больше внимания лечению заболевания. Автоматизированная система скрининга могла бы также обеспечить большую точность диагностики, исключив из процесса субъективный человеческий фактор.

В этом примере мы исследуем полезность машинного обучения для выявления рака, применив алгоритм k-NN к исследованиям клеток, полученных при биопсии, у женщин с аномальными новообразованиями молочной железы.

Этап 1. Сбор данных

Воспользуемся набором данных *Breast Cancer Wisconsin (Diagnostic)*, полученным из репозитория машинного обучения UCI, расположенного по адресу <http://archive.ics.uci.edu/ml>. Эти данные были переданы в репозиторий исследователями из Висконсинского университета и включают в себя измерения, полученные по оцифрованным изображениям материала, извлеченного тонкой иглой из образований в области молочной железы. Это характеристики клеточных ядер, обнаруженных на цифровом изображении.



Подробнее об этом наборе данных вы можете прочитать в статье: Mangasarian O.L., Street W.N., Wolberg W.H. Breast Cancer Diagnosis and Prognosis via Linear Programming. *Operations Research*, 1995. Vol. 43. P. 570–577.

Данные о раке молочной железы включают в себя 569 примеров биопсии рака, каждый из которых имеет 32 признака: идентификационный номер, вид рака и еще 30 числовых лабораторных измерений. Диагноз обозначается как "М" для злокачественного (malignant) или "В" для доброкачественного (benign) образования.

Тридцать числовых измерений включают в себя среднее арифметическое, стандартную ошибку и наилучшее (то есть наибольшее) значение десяти различных характеристик оцифрованных ядер клеток. К ним относятся:

- радиус;
- текстура;
- периметр;
- площадь;
- гладкость;
- компактность;
- вогнутость;
- вогнутые точки;
- симметричность;
- фрактальные размеры.

На основании перечисленных характеристик можно сделать вывод, что все признаки, похоже, описывают форму и размер клеточных ядер. Если вы не онколог, то вряд ли узнаете, как именно каждый из них относится к

доброкачественным или злокачественным образованиям. Эти паттерны будут раскрыты, когда мы продолжим процесс машинного обучения.

Этап 2. Исследование и подготовка данных

Изучим данные и посмотрим, сможем ли мы выявить взаимосвязь между ними. Подготовим данные для использования с помощью метода обучения k-NN.



Если вы планируете продолжить выполнение примера, сохраните файл `wisc_bc_data.csv` в своем рабочем каталоге R (файл вы можете скачать вместе с остальными материалами по адресу <https://github.com/dataspelunking/MLwR/>). Для этой книги набор данных был немного изменен по сравнению с его исходной формой. В частности, добавлена строка заголовка, а строки данных расположены случайным образом.

Для начала импортируем файл данных в формате CSV, как мы это делали в предыдущих главах, и сохраним данные о раке молочной железы из Висконсинского университета во фрейме данных `wbcd`:

```
> wbcd <- read.csv("wisc_bc_data.csv",  
stringsAsFactors = FALSE)
```

Воспользовавшись командой `str(wbcd)`, мы можем подтвердить, что данные, как и ожидалось, структурированы и представляют собой 569 примеров с 32 признаками. Вот несколько первых строк результатов этой функции:

```
'data.frame': 569 obs. of 32 variables:$  
id : int 87139402 8910251  
905520 ...$ diagnosis : chr "B" "B"  
"B" "B" ...$ radius_mean : num 12.3 10.6  
11 11.3 15.2 ...$ texture_mean : num 12.4  
18.9 16.8 13.4 13.2 ...$ perimeter_mean :  
num 78.8 69.3 70.9 73 97.7 ...$  
area_mean : num 464 346 373 385 712  
...
```

Первая переменная представляет собой целочисленную переменную с именем `id`. Поскольку это просто *уникальный идентификатор (ID)* пациента, то он не предоставляет полезной информации и нам нужно будет исключить его из модели.



Независимо от метода машинного обучения ID-переменные всегда должны быть исключены. Пренебрежение этим может привести к ошибочным выводам, поскольку идентификатор можно использовать для правильного прогнозирования каждого примера. Следовательно, модель, включающая столбец идентификатора, почти наверняка пострадает от переобучения и будет плохо обобщена для будущих данных.

Давайте вообще отбросим признак `id`. Поскольку он расположен в первом столбце, то, чтобы его исключить, нужно сделать копию фрейма данных `wbcd` без столбца 1:

```
> wbcd <- wbcd[-1]
```

Следующая переменная, `diagnosis`, представляет особый интерес, так как это результат, который мы надеемся спрогнозировать. Этот признак указывает на то, относится ли данный образец к доброкачественному или злокачественному новообразованию. Результат выполнения функции `table()` показывает, что 357 новообразований являются доброкачественными, а 212 — злокачественными:

```
> table(wbcd$diagnosis)  B      M357  212
```

Многие классификаторы машинного обучения на языке R требуют, чтобы целевая функция была закодирована как фактор, поэтому нам нужно будет перекодировать переменную `diagnosis`. Мы также воспользуемся этой возможностью, чтобы присвоить значениям "B" и "M" более информативные метки, используя параметр `labels`:

```
> wbcd$diagnosis <- factor(wbcd$diagnosis,
levels = c("B", "M"), labels = c("Benign",
"Malignant"))
```

Глядя на результаты `prop.table()`, мы увидим, что значения помечены как `Benign` и `Malignant` и равны 62,7 и 37,3 % всех новообразований соответственно:

```
> round(prop.table(table(wbcd$diagnosis)) *
100, digits =
1)  Benign      Malignant      62.7      37.3
```

Остальные 30 признаков являются числовыми и, как и ожидалось, состоят из трех разных измерений десяти характеристик. Мы рассмотрим только три таких признака:

```
> summary(wbcd[c("radius_mean", "area_mean",
"smoothness_mean")])  radius_mean      area_mean
smoothness_meanMin.      : 6.981  Min.      :
143.5  Min.      :0.052631st Qu.  :11.700  1st Qu.  :
420.3  1st Qu.  :0.08637Median    :13.370  Median    :
551.1  Median    :0.09587Mean      :14.127  Mean      :
```

```
654.9    Mean      :0.096363rd Qu.  :15.780   3rd Qu.  :
782.7    3rd
Qu. :0.10530Max.      :28.110   Max.      :2501.0   Max
.      :0.16340
```

Замечаете ли вы что-нибудь проблематичное в значениях, когда смотрите на эти три признака, расположенных рядом? Напомню, что вычисление расстояния для k-NN очень зависит от масштаба измерения исходных признаков. Поскольку гладкость (`smoothness`) колеблется от 0,05 до 0,16, а площадь (`area`) — от 143,5 до 2501,0, то при вычислении расстояния площадь будет влиять намного больше, чем гладкость. Это может вызвать у нашего классификатора проблемы, поэтому применим нормализацию, чтобы изменить масштаб функций до стандартного диапазона значений.

Преобразование: нормализация числовых данных

Для того чтобы нормализовать эти функции, нам нужно создать в R функцию `normalize()`. Эта функция принимает вектор числовых значений `x`, вычитает из каждого значения `x` минимальное значение `x`, делит его на диапазон значений `x` и, наконец, возвращает результирующий вектор. Код функции выглядит следующим образом:

```
> normalize <- function(x) { return ((x -
min(x)) / (max(x) - min(x))) }
```

После выполнения предыдущего кода функция `normalize()` доступна для использования в R. Проверим эту функцию на нескольких векторах:

```
> normalize(c(1, 2, 3, 4, 5))[1] 0.00 0.25 0.50
0.75 1.00> normalize(c(10, 20, 30, 40, 50))[1]
0.00 0.25 0.50 0.75 1.00
```

Функция работает правильно. Несмотря на то что значения во втором векторе в десять раз больше, чем в первом, после нормализации оба вектора выглядят абсолютно одинаково.

Теперь мы можем применить функцию `normalize()` к числовым объектам фрейма данных. Вместо того чтобы нормализовать каждую из 30 числовых переменных по отдельности, воспользуемся одной из функций R для автоматизации процесса.

Функция `lapply()` принимает список и применяет функцию нормализации к каждому элементу этого списка. Поскольку фрейм данных представляет собой список векторов равной длины, мы можем использовать функцию `lapply()`, чтобы применить `normalize()` к каждому признаку фрейма данных. На последнем этапе список,

возвращаемый `lapply()`, преобразуется во фрейм данных с помощью функции `as.data.frame()`. Полностью процесс выглядит так:

```
> wbcd_n <- as.data.frame(lapply(wbcd[2:31],  
normalize))
```

Проще говоря, эта команда применяет функцию `normalize()` к столбцам с 2 по 31 во фрейме данных `wbcd`, преобразует результат во фрейм данных и присваивает ему имя `wbcd_n`. Суффикс `_n` здесь используется как напоминание о том, что значения в `wbcd` теперь нормализованы.

Чтобы убедиться, что преобразование выполнено правильно, посмотрим на сводные данные по одной переменной:

```
> summary(wbcd_n$area_mean)Min.      1st  
Qu.    Median      Mean      3rd  
Qu.    Max.0.0000    0.1174    0.1729    0.2169  
0.2711    1.0000
```

Как и ожидалось, переменная `area_mean`, которая изначально составляла от 143,5 до 2501,0, теперь варьируется от 0 до 1.

Подготовка данных: создание тренировочных и тестовых наборов данных

Несмотря на то что все 569 результатов биопсии помечены статусом «доброкачественный» или «злокачественный», прогнозировать то, что мы и так уже знаем, едва ли интересно. Кроме того, любые показатели эффективности, которые мы получим во время обучения, могут вводить в заблуждение, так как мы не знаем степень аппроксимации данных или того, насколько хорошо обучаемый будет делать обобщения для новых случаев. Поэтому более интересно то, насколько хорошо обучаемый работает с данными, не входящими в тренировочный набор. Если бы у нас был доступ к лаборатории, то мы могли бы предоставить обучаемому результаты биопсии, полученных из следующих 100 новообразований с неизвестным статусом рака, и увидеть, насколько хорошо прогнозы обучающейся машины совпадают с диагнозами, полученными традиционными методами.

В отсутствие таких данных мы можем смоделировать подобный сценарий, разделив данные на две части: тренировочный набор данных, который будет использоваться для построения модели `k-NN`, и тестовый набор данных для оценки точности прогнозирования модели. Мы будем использовать первые 469 записей для тренировочного набора данных и оставшиеся 100 — для имитации новых пациентов.

Используя методы извлечения данных, представленные в главе 2, мы разделим фрейм данных `wbcd_n` на фреймы `wbcd_train` и `wbcd_test`:

```
> wbcd_train <- wbcd_n[1:469, ] > wbcd_test <-  
wbcd_n[470:569, ]
```

Если эти команды сбивают вас с толку, вспомните, что данные извлекаются из фреймов данных с использованием синтаксиса `[row, col]`. Пустое значение на месте строки или столбца указывает на то, что здесь должны быть включены все строки или столбцы. Следовательно, первая строка кода запрашивает строки с 1 по 469 со всеми столбцами, а вторая — 100 строк с 470 по 569, также со всеми столбцами.



При построении тренировочного и тестового наборов данных важно, чтобы каждый набор данных представлял собой репрезентативное подмножество полного набора. Записи `wbcd` уже были перемешаны случайным образом, поэтому мы могли просто извлечь 100 последовательных записей, чтобы создать тестовый набор данных. Если бы данные были упорядочены в хронологическом порядке или по группам схожих значений, это было бы неуместно. В таких случаях потребуются методы случайной выборки. Случайная выборка будет рассмотрена в главе 5.

После того как мы сформировали нормализованные тренировочный и тестовый наборы данных, нам нужно исключить из них целевую переменную — `diagnosis`. Для обучения модели `k-NN` эти метки классов нужно хранить в факторных векторах, разделенных на тренировочный и тестовый наборы данных:

```
> wbcd_train_labels <- wbcd[1:469, 1] >  
wbcd_test_labels <- wbcd[470:569, 1]
```

Этот код принимает фактор `diagnosis` из первого столбца таблицы данных `wbcd` и создает векторы `wbcd_train_labels` и `wbcd_test_labels`. Мы будем использовать их на следующих этапах для обучения и оценки нашего классификатора.

Шаг 3. Обучение модели на данных

Благодаря тренировочным данным и вектору меток мы теперь готовы классифицировать тестовые записи. Для алгоритма `k-NN` этап обучения фактически не включает в себя построение модели; процесс обучения ленивого обучаемого, такого как `k-NN`, подразумевает только сохранение входных данных в структурированном формате.

Для классификации наших тестовых экземпляров мы воспользуемся реализацией k-NN из пакета `class`. Эта реализация представляет собой набор базовых R-функций классификации. Если указанный пакет еще не установлен в вашей системе, то установите его, введя следующую команду:

```
> install.packages("class")
```

Для того чтобы загрузить пакет во время сеанса, в течение которого вы хотите использовать функции этого пакета, просто введите команду `library(class)`.

Функция `knn()` из пакета `class` представляет собой стандартную, классическую реализацию алгоритма k-NN. Для каждого экземпляра тестовых данных эта функция будет идентифицировать k ближайших соседей, используя евклидово расстояние, где k — число, заданное пользователем. Тестовый экземпляр классифицируется путем «голосования» среди k ближайших соседей — другими словами, это означает присвоение экземпляру того же класса, что и у большинства соседей. В случае равенства голосов решение принимается случайным образом.



В других R-пакетах есть несколько иных функций k-NN, в которых реализованы более сложные или более эффективные варианты алгоритма. Если вы столкнетесь с ограничениями, используя `knn()`, то поищите алгоритм k-NN в Comprehensive R Archive Network (CRAN).

Обучение и классификация с использованием функции `knn()` выполняются с помощью одной команды, имеющей четыре параметра.

Синтаксис классификации kNN

Использование функции `knn()` из пакета `class`

Построение классификатора и создание прогнозов:

```
p <- knn(train, test, class, k)
```

`train` – это фрейм данных, содержащий тренировочные данные в цифровом виде;

`test` – это фрейм данных, содержащий тестовые данные в цифровом виде;

`class` – это факторный вектор, содержащий значения класса для каждой строки тренировочных данных;

`k` – это целое число, означающее количество ближайших соседей.

Функция возвращает факторный вектор спрогнозированных классов для каждой строки тестового фрейма данных.

Пример:

```
wbcd_pred <- knn(train = wbcd_train, test=wbcd_test,  
                cl = wbcd_train_labels, k = 3)
```

Теперь у нас есть почти все, что нужно, чтобы применить алгоритм k-NN к этим данным. Мы разделили данные на тренировочный и тестовый наборы с одинаковыми числовыми характеристиками. Метки для тренировочных данных хранятся в отдельном факторном векторе. Осталось указать только параметр `k`, определяющий количество соседей, которых необходимо включить в голосование.

Поскольку тренировочные данные включают в себя 469 экземпляров, мы могли бы попробовать `k=21` — нечетное число, примерно равное квадратному корню из 469. Если алгоритму придется выбирать из двух категорий, то использование нечетного числа исключит возможность разделения голосов на две равные группы.

Теперь можно воспользоваться функцией `knn()` для классификации тестовых данных:

```
> wbcd_test_pred <- knn(train = wbcd_train,  
                        test = wbcd_test,                          cl =  
                        wbcd_train_labels, k = 21)
```

Функция `knn()` возвращает факторный вектор спрогнозированных меток для всех объектов из набора данных `wbcd_test`. Мы назвали эти прогнозы `wbcd_test_pred`.

Шаг 4. Оценка эффективности модели

Следующий шаг — оценка того, насколько хорошо спрогнозированные классы, записанные в векторе `wbcd_test_pred`, соответствуют действительным значениям в векторе `wbcd_test_labels`. Для этого можно воспользоваться функцией `CrossTable()` из пакета `gmodels`,

который был описан в главе 2. Если вы еще не установили его, сделайте это с помощью команды `install.packages("gmodels")`.

Загрузив пакет с помощью команды `library(gmodels)`, можно построить перекрестную таблицу, указывающую на соответствие между векторами спрогнозированных и действительных меток.

Аргумент `prop.chisq=FALSE` позволяет удалить из вывода ненужные значения хи-квадрат:

```
> CrossTable(x = wbcd_test_labels, y =  
wbcd_test_pred, prop.chisq = FALSE)
```

Результирующая таблица выглядит так:

wbcd_test_labels	wbcd_test_pred		Row Total
	Benign	Malignant	
Benign	61 1.000 0.968 0.610	0 0.000 0.000 0.000	61 0.610
Malignant	2 0.051 0.032 0.020	37 0.949 1.000 0.370	39 0.390
Column Total	63 0.630	37 0.370	100

Проценты ячеек указывают на долю значений, которые делятся на четыре категории. В верхней левой ячейке указаны *правильно спрогнозированные отрицательные* результаты; 61 из 100 значений — это те случаи, когда новообразование было доброкачественным и алгоритм k-NN правильно его идентифицировал как таковое. В правой нижней ячейке указаны *правильно спрогнозированные положительные* результаты, то есть те случаи, когда классификатор и клинически определенная метка согласны с тем, что данное новообразование является злокачественным. В общей сложности 37 из 100 прогнозов были правильно распознаны как положительные.

Ячейки, расположенные в другой диагонали, содержат количество примеров, для которых прогноз k-NN не соответствует истинной метке. Два примера в нижней левой ячейке являются *ложными отрицательными* результатами; в этом случае прогноз дал доброкачественное значение, но на самом деле опухоль была злокачественной. Ошибки в этом в этом случае могут иметь серьезные последствия, поскольку пациент будет думать, что у него нет рака, тогда как на самом деле болезнь прогрессирует.

Верхняя правая ячейка должна содержать ложные положительные результаты, если таковые есть. Такие значения возникают, когда модель классифицирует новообразование как злокачественное, а в действительности оно является доброкачественным. Подобные ошибки менее опасны, чем ложный отрицательный результат, однако их также

следует избегать, поскольку они могут привести к дополнительным финансовым затратам системы здравоохранения и к излишнему стрессу для пациента, дополнительным анализам или ненужному лечению.



При желании можно было бы полностью устранить ложные отрицательные классификации, отмечая каждое новообразование как злокачественное. Однако очевидно, что это нереальная стратегия. Тем не менее она показывает, что прогнозирование предполагает достижение баланса между ложными положительными и ложными отрицательными классификациями. В главе 10 мы рассмотрим методы оценки точности прогнозирования, которые можно использовать для оптимизации эффективности алгоритма с учетом стоимости каждого типа ошибок.

В общей сложности метод k -NN неправильно классифицировал 2 из 100, или 2 % новообразований. Такая 98%-ная точность кажется впечатляющей с учетом всего нескольких строк R-кода, однако стоит попробовать еще немного улучшить модель, чтобы посмотреть, сможем ли мы повысить эффективность и сократить количество ошибочно классифицированных значений, особенно с учетом опасности ложных отрицательных результатов.

Шаг 5. Повышение эффективности модели

Проверим два простых варианта предыдущего классификатора. Сначала применим альтернативный метод для масштабирования числовых функций. Затем попробуем несколько разных значений k .

Преобразование: стандартизация по z-оценке

Обычно при классификации k -NN используется нормализация, однако для масштабирования признаков в наборе данных о раковых новообразованиях более подходящий способ — стандартизация по z-оценке. Поскольку значения, стандартизированные по z-оценке, не имеют заранее определенного минимума и максимума, экстремальные значения не сжимаются к центру. Даже не имея медицинского образования, можно предположить, что злокачественная опухоль может привести к экстремальным выбросам, поскольку опухоли имеют свойство неконтролируемо расти. С учетом этого, возможно, было бы разумнее допустить учет выбросов при расчете расстояния. Посмотрим, улучшит ли стандартизация по z-оценке точность прогнозирования.

Чтобы стандартизировать вектор, можно задействовать встроенную R-функцию `scale()`, которая по умолчанию масштабирует значения, используя стандартизацию по z-оценке. Функция `scale()` может быть применена непосредственно к фрейму данных, поэтому нет необходимости

использовать функцию `lapply()`. Для того чтобы создать версию данных `wbcd`, стандартизированную по z-оценке, можно воспользоваться следующей командой:

```
> wbcd_z <- as.data.frame(scale(wbcd[-1]))
```

Эта команда изменяет масштаб всех признаков, за исключением `diagnosis`, значения которого находятся в первом столбце, и сохраняет результат в виде фрейма данных `wbcd_z`.

Суффикс `_z` говорит о том, что значения были преобразованы по z-оценке.

Для того чтобы убедиться, что преобразование было выполнено правильно, посмотрим на сводные данные:

```
> summary(wbcd_z$area_mean)
  Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
1.4530 -0.6666      -0.6666      -0.6666      0.3632      5.2460
```

Среднее арифметическое переменной, стандартизированной по z-оценке, всегда должно быть равно нулю, а диапазон должен быть достаточно компактным. Если z-оценка больше 3 или меньше -3, это указывает на чрезвычайно редкое значение. С учетом этого, изучая сводные данные, можно сделать вывод, что преобразование, похоже, выполнено правильно.

Теперь нужно разделить данные, преобразованные по z-оценке, на тренировочный и тестовый наборы и классифицировать тестовые экземпляры с помощью функции `knn()`. Затем сравним спрогнозированные метки с реальными, используя функцию `CrossTable()`:

```
> wbcd_train <- wbcd_z[1:469, ]
> wbcd_test <- wbcd_z[470:569, ]
> wbcd_train_labels <- wbcd[1:469, 1]
> wbcd_test_labels <- wbcd[470:569, 1]
> wbcd_test_pred <- knn(train = wbcd_train,
  test = wbcd_test,
  cl = wbcd_train_labels, k = 21)
> CrossTable(x = wbcd_test_labels, y = wbcd_test_pred,
  prop.chisq = FALSE)
```

К сожалению, как видно ниже, результаты нового преобразования показывают даже небольшое снижение точности. В тех случаях, где в прошлый раз было правильно классифицировано 98 % примеров, теперь правильно классифицировано лишь 95 %. Что еще хуже, не получилось лучше классифицировать опасные ложные отрицательные результаты.

wbc_d_test_labels	wbc_d_test_pred		Row Total
	Benign	Malignant	
Benign	61	0	61
	1.000	0.000	0.610
	0.924	0.000	
	0.610	0.000	
Malignant	5	34	39
	0.128	0.872	0.390
	0.076	1.000	
	0.050	0.340	
Column Total	66	34	100
	0.660	0.340	

Тестирование альтернативных значений k

Можно оптимизировать эффективность модели k -NN, изучив ее возможности при различных значениях k . Используя нормализованные тренировочный и тестовый наборы данных, мы классифицировали те же 100 записей с использованием нескольких вариантов k . В табл. 3.4 показано количество ложных отрицательных и ложных положительных результатов для каждой итерации.

Классификатор ни разу не сработал идеально, однако подход 1-NN позволил избежать некоторых ложных отрицательных результатов — правда, при этом увеличилось количество ложных положительных результатов. Важно помнить, однако, что неразумно слишком точно адаптировать этот алгоритм к тестовым данным; ведь другой набор из 100 записей о пациентах, скорее всего, будет отличаться от тех, которые использованы для измерения эффективности алгоритма.

Таблица 3.4

Значение k	Ложные отрицательные результаты	Ложные положительные результаты	Некорректно распознанные результаты, %
1	1	3	4
5	2	0	2
11	3	0	3
15	3	0	3
21	2	0	2
27	4	0	4



Если необходимо гарантировать, что в дальнейшем алгоритм обучения будет обобщать данные, можно создать несколько наборов из 100 пациентов. Затем расположить наборы в случайном порядке и повторно протестировать результат. Такие методы тщательной оценки эффективности моделей машинного обучения будут рассмотрены в главе 10.

Резюме

В этой главе мы рассмотрели алгоритмы классификации k -NN. В отличие от многих других алгоритмов классификации метод *ближайших соседей* не обучается. Он просто в точности сохраняет тренировочные данные. Затем непроверенные тестовые примеры сопоставляются с наиболее похожими записями тренировочного набора с помощью функции вычисления расстояния и немаркированному примеру присваивается метка его соседей.

Хотя k -NN очень простой алгоритм, он способен решать чрезвычайно сложные задачи, такие как выявление раковых новообразований. С использованием всего нескольких простых строк R-кода в 98 % случаев было правильно установлено, является новообразование злокачественным или доброкачественным.

В следующей главе мы разберем метод классификации для оценки вероятности того, что наблюдение попадает в определенные категории. Интересно сравнить, насколько этот подход отличается от k -NN. В главе 9 мы поговорим о методе, очень похожем на k -NN, но меры расстояния в нем используются для совершенно другой задачи обучения.

4. Вероятностное обучение: классификация с использованием наивного байесовского классификатора

Когда метеоролог составляет прогноз погоды, осадки в нем, как правило, описываются фразами наподобие такой: «вероятность дождя — 70 %».

Такие прогнозы называются отчетами о вероятности осадков.

Задумывались ли вы, как рассчитываются эти прогнозы? Вопрос непростой, потому что на самом деле неизвестно, пойдет дождь или нет.

Прогнозы погоды основаны на вероятностных методах, которые направлены на описание неопределенности. Это методы экстраполяции будущих событий, основанные на данных о прошлых событиях. В случае с погодой вероятность дождя описывается на основании того, какая доля дней с аналогичными атмосферными условиями в прошлом сопровождалась выпадением осадков. Вероятность дождя в 70 % означает, что в семи из десяти прошлых случаев с аналогичными условиями наблюдалось выпадение осадков.

В этой главе мы познакомимся с наивным байесовским классификатором. В алгоритме вероятности используются практически так же, как и при прогнозе погоды. Рассмотрим следующие темы:

- основные принципы вероятности;
- специализированные методы и структуры данных, необходимые для анализа текстовых данных с помощью R;

- как использовать наивный байесовский классификатор для создания фильтра нежелательных СМС.

Если ранее вы изучали статистику, то некоторые материалы этой главы можно бегло просмотреть. Однако полезно освежить знания о теории вероятностей, так как основные принципы объясняют, почему наивный байесовский классификатор получил такое интересное название.

Наивный байесовский классификатор

Основные положения статистики, необходимые для понимания алгоритма наивной байесовской классификации, известны на протяжении нескольких веков. Начало этой технологии положил математик XVII века Томас Байес, который разработал основополагающие принципы вероятности наступления событий и того, как следует пересматривать вероятности в свете поступления дополнительной информации. Эти принципы легли в основу так называемых *байесовских методов*.

Позже мы рассмотрим эти методы более подробно. Пока достаточно сказать, что вероятность — это число в диапазоне от 0 до 1 (то есть от 0 до 100 %), которое отражает вероятность того, что событие произойдет, с учетом имеющихся данных. Чем ниже вероятность, тем менее вероятно наступление события. Нулевая вероятность указывает на то, что событие точно не произойдет, в то время как вероятность, равная 1, указывает на то, что событие произойдет абсолютно точно.

Классификаторы, основанные на байесовских методах, используют тренировочные данные для расчета вероятности каждого результата на базе данных, представленных в виде значений признаков. Когда классификатор применяется к неклассифицированным данным, он использует эти вычисленные вероятности, чтобы спрогнозировать наиболее вероятный класс для нового примера. Идея проста, но она ведет к методу, результаты которого сравнимы с более сложными алгоритмами. На практике байесовские алгоритмы использовались для решения таких задач, как:

- классификация текста, например фильтрация нежелательной почты (спама);
- обнаружение вторжений и аномалий в компьютерных сетях;
- диагностика заболеваний по совокупности наблюдаемых симптомов.

Как правило, байесовские алгоритмы лучше всего применимы к тем задачам, в которых необходимо одновременно учесть информацию, состоящую из многочисленных атрибутов, чтобы оценить суммарную вероятность результата. Многие алгоритмы машинного обучения игнорируют слабо выраженные особенности, однако байесовские методы используют все имеющиеся доказательства, чтобы слегка изменить прогнозы. Отсюда следует, что, даже если большое количество объектов

оказывает относительно незначительное воздействие, их совокупное влияние в байесовской модели может оказаться довольно большим.

Основные понятия байесовских методов

Прежде чем перейти к наивному байесовскому алгоритму, стоит уделить некоторое время определению понятий, которые используются в байесовских методах. Теория вероятностей Байеса основана на идее, что предполагаемая вероятность *события* или потенциальный исход, должны основываться на уже имеющихся данных, полученных в результате нескольких *исходов*, или на возможностях возникновения события.

В табл. 4.1 показаны события и исходы для нескольких реальных результатов.

Таблица 4.1

Событие	Исход
Подбрасывание монеты	Орел
Один день	Дождливая погода
Получено почтовое сообщение	Сообщение является спамом
Президентские выборы	Кандидат выбран президентом
Покупка лотерейного билета	Выигрыш в лотерею

Байесовские методы дают представление о том, как можно оценить вероятность наступления этих событий по наблюдаемым данным. Для того чтобы понять, как это происходит, нам нужно разобраться в том, что такое вероятность.

Что такое вероятность

Вероятность наступления события оценивается по наблюдаемым данным путем деления количества исходов, в которых произошло событие, на общее количество исходов. Например, если дождь шел три дня из десяти с такими же условиями, как сегодня, то вероятность дождя сегодня можно оценить как $3 / 10 = 0,30$, или 30 %. Аналогично если 10 из 50 предыдущих сообщений электронной почты были спамом, то вероятность того, что любое входящее сообщение является спамом, можно оценить как $10 / 50 = 0,20$, или 20 %.

Чтобы обозначить эти вероятности, воспользуемся обозначением $P(A)$, что означает вероятность наступления события A . Например, $P(\text{дождь}) = 0,30$, а $P(\text{спам}) = 0,20$.

Вероятность всех возможных исходов всегда должна быть равна 1, потому что любой исход всегда приводит к тому, что происходит какой-то результат. Таким образом, если в исходе есть два результата, которые не могут произойти одновременно, — например, дождливый и солнечный день или спам и полезное письмо (не спам), — то, зная вероятность одного из этих результатов, мы будем знать вероятность другого. Например,

учитывая, что $P(\text{спам}) = 0,20$, можно вычислить $P(\text{не спам}) = 1 - 0,20 = 0,80$. Этот подход работает, поскольку спам и не спам являются взаимоисключающими и исчерпывающими событиями — другими словами, они не могут происходить одновременно и являются единственно возможными результатами.

Далее, поскольку событие не может происходить и не происходить одновременно, то событие и его *дополнение* — событие, содержащее результаты, при которых интересующее событие не происходит, — всегда являются взаимоисключающими и исчерпывающими. Дополнение к событию A обычно обозначается как A^c . Кроме того, для обозначения вероятности того, что событие A не произойдет, иногда используется сокращенная запись $P(\neg A)$. Например, обозначение $P(\neg \text{спам}) = 0,80$ эквивалентно $P(A^c)$.

Чтобы лучше понять суть события и их дополнения, представьте себе двумерное пространство, которое разделено на вероятности для каждого события. На рис. 4.1 прямоугольник означает все возможные результаты для сообщения электронной почты. Круг соответствует 20%-ной вероятности того, что сообщение является спамом. Оставшиеся 80 % представляют собой дополнение $P(\neg \text{спам})$, то есть вероятность того, что сообщение не является спамом.

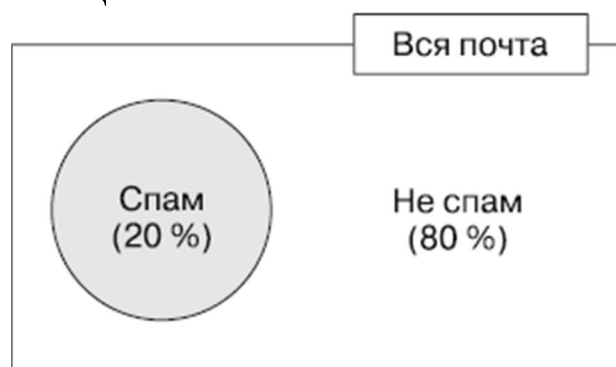


Рис. 4.1. Пространство вероятности всей электронной почты можно представить как разделы со спамом и не спамом

Суммарная вероятность

Часто мы заинтересованы в мониторинге нескольких невзаимоисключающих событий для одного и того же исхода. Если определенные события происходят одновременно с интересующим нас событием, то их можно использовать для прогнозирования. Рассмотрим, например, второе событие, основанное на результате, когда в сообщении электронной почты содержится слово «виагра». Приведенная выше диаграмма с учетом второго события может выглядеть так, как показано на рис. 4.2.

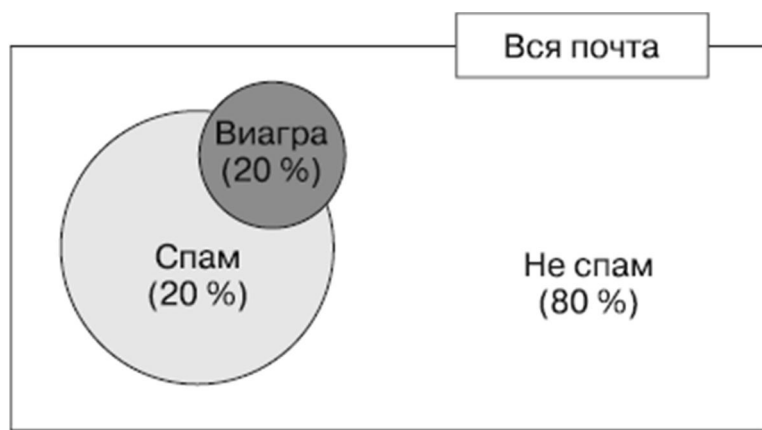


Рис. 4.2. Невзаимоисключающие события изображаются как перекрывающиеся области

Обратите внимание, что на рис. 4.2 круг «Виагра» не полностью заполняет круг «Спам» и не полностью ограничен этим кругом. Это означает, что не все спамовые сообщения содержат слово «виагра» и не любое письмо со словом «виагра» является спамом. Тем не менее, поскольку это слово очень редко встречается вне спама, его наличие во входящем сообщении будет подтверждением того, что это сообщение является спамом.

Чтобы увеличить масштаб изображения и подробнее рассмотреть эти перекрывающиеся круги, воспользуемся визуализацией, известной как *диаграмма Венна*. Ее впервые использовал в конце XIX века математик Джон Венн. На ней с помощью кругов показано перекрытие множеств. В большинстве диаграмм Венна, и здесь в том числе, диаметр кругов и степень перекрытия не имеют значения.

Он лишь служит напоминанием о распределении вероятности для всех комбинаций событий. Диаграмма Венна для спама и виагры может быть изображена как на рис. 4.3.

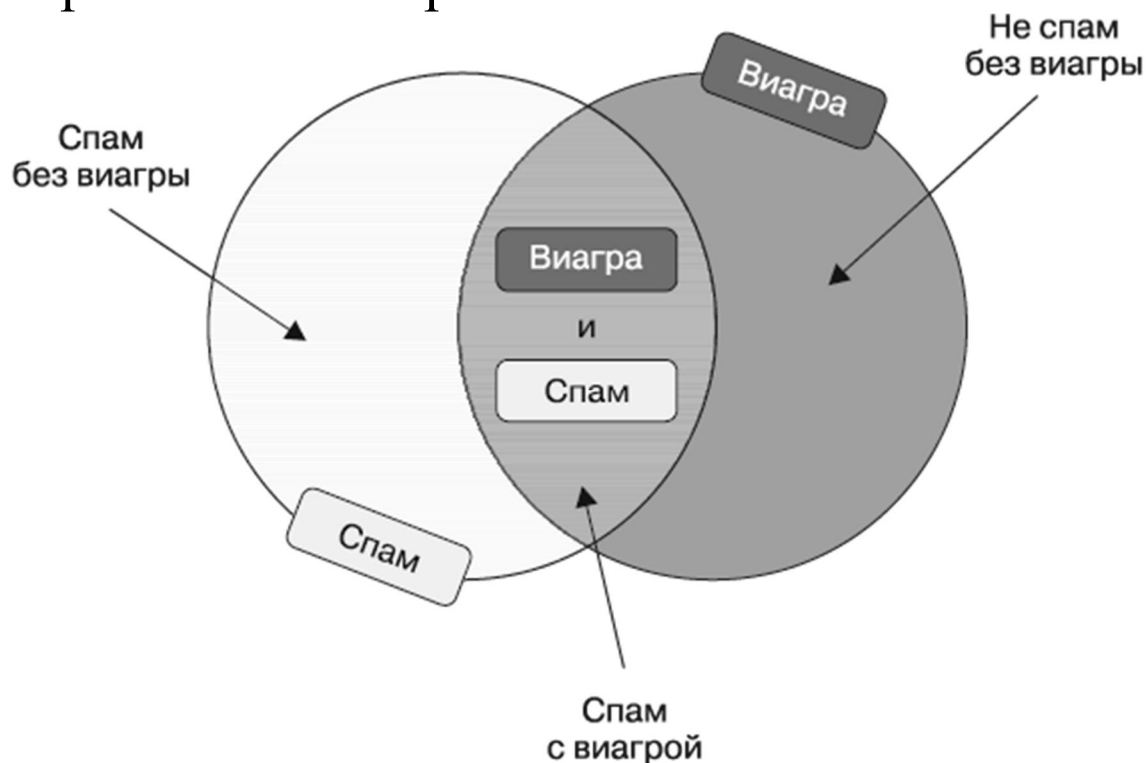


Рис. 4.3. На диаграмме Венна видно совпадение событий спама и наличия в сообщениях слова «виагра»

Нам известно, что 20 % всех сообщений были спамом (левый кружок), а 5 % всех сообщений содержали слово «виагра» (правый кружок). Мы хотели бы количественно определить степень совпадения между этими двумя пропорциями. Другими словами, мы надеемся оценить вероятность

совпадения $P(\text{спам})$ и $P(\text{виагра})$, что можно записать как $P(\text{спам} \cap \text{виагра})$. Перевернутый символ \cup означает пересечение двух событий; обозначение $A \cap B$ означает событие, при котором происходит как A , так и B .

Вычисление $P(\text{спам} \cap \text{виагра})$ зависит от *суммарной вероятности* двух событий, или, другими словами, от того, как вероятность одного события связана с вероятностью другого. Если эти события совершенно не связаны, то они называются *независимыми событиями*. Это не означает, что независимые события не могут происходить одновременно; независимость события лишь подразумевает, что знание результата одного события не дает никакой информации о результате другого. Например, результат того, что при подбрасывании монеты выпадет орел, не зависит от того, какая в этот день погода — дождливая или солнечная.

Если бы все события были независимыми, то было бы невозможно спрогнозировать одно событие, наблюдая другое. Иначе говоря, *зависимые события* являются основой прогнозного моделирования (рис. 4.4). Подобно тому как облака предвещают дождливый день, наличие в сообщении слова «виагра» говорит о том, что, скорее всего, это письмо является спамом.

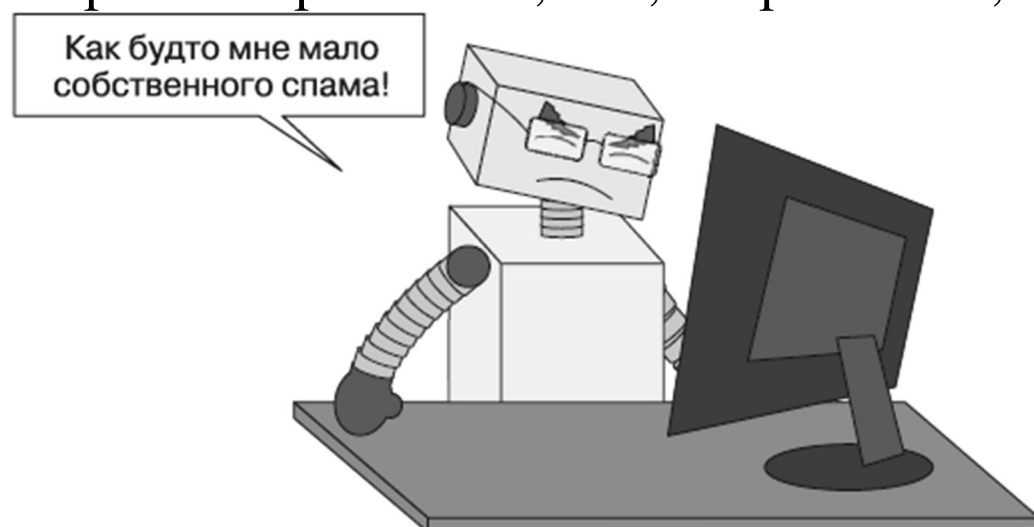


Рис. 4.4. Зависимые события нужны машине, чтобы научиться выявлять полезные закономерности

Вычисление вероятности зависимых событий немного сложнее, чем в случае с независимыми событиями. Если бы $P(\text{спам})$ и $P(\text{виагра})$ были независимыми, то мы могли бы легко рассчитать $P(\text{спам} \cap \text{виагра})$ — вероятность того, что оба эти события произошли одновременно. Поскольку 20 % всех сообщений являются спамом, а 5 % всех электронных писем содержат слово «виагра», то можно предположить, что 1 % всех сообщений является спамом, в котором содержится слово «виагра», потому что $0,05 \times 0,20 = 0,01$. В более общем случае для независимых событий A и B вероятность их наступления можно вычислить так: $P(A \cap B) = P(A) \times P(B)$.

Однако $P(\text{спам})$ и $P(\text{виагра})$, очевидно, сильно зависят друг от друга, а это означает, что такой расчет неверен. Чтобы получить надлежащую оценку, нам нужно использовать более точную формулировку взаимосвязи между этими двумя событиями, которая основана на углубленных байесовских методах.

Вычисление условной вероятности по теореме Байеса

Взаимосвязь между зависимыми событиями может быть описана с помощью *теоремы Байеса*, которая позволяет пересмотреть оценку вероятности наступления одного события с учетом сведений, предоставленных другим событием. Одна из формулировок этой теоремы выглядит следующим образом:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} .$$

Обозначение $P(A | B)$ читается как «вероятность наступления события A с учетом того, что произошло событие B ». Это называется *условной вероятностью*, поскольку событие A зависит (то есть выполняется при условии) от того, что событие B уже произошло. Согласно теореме Байеса оценка $P(A | B)$ должна основываться на $P(A \cap B)$ — мере того, как часто события A и B происходят вместе, и на $P(B)$ — мере того, как часто вообще происходит событие B .

Теорема Байеса утверждает, что наилучшая оценка $P(A | B)$ — это отношение исходов, в которых событие A произошло вместе с B , ко всем исходам, когда произошло событие B . Это подразумевает, что вероятность события A выше, если A и B часто происходят вместе всякий раз, когда происходит событие B . Обратите внимание, что эта формула корректирует $P(A \cap B)$ для вероятности наступления события B . Если B происходит очень редко, то $P(B)$ и $P(A \cap B)$ всегда будут маленькими; однако если события A и B почти всегда происходят вместе, то $P(A | B)$ будет высоким независимо от вероятности наступления события B .

По определению $P(A \cap B) = P(A | B) \times P(B)$ — вывод, который легко получить, применив законы алгебры к предыдущей формуле. Перегруппировав эту формулу еще раз, при этом зная, что $P(A \cap B) = P(B \cap A)$, получим результат $P(A \cap B) = P(B | A) \times P(A)$, который затем можно использовать в следующей формулировке теоремы Байеса:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)} .$$

В сущности, это традиционная формулировка теоремы Байеса. Когда мы применим ее к машинному обучению, станет ясно, почему это так. Но пока, чтобы лучше понять, как работает теорема Байеса на практике, вернемся к гипотетическому фильтру спама.

Если не знать содержимое входящего сообщения, наилучшей оценкой его статуса — является ли сообщение спамом — будет $P(\text{спам})$, вероятность того, что любое предыдущее сообщение было спамом. Эта оценка называется *априорной вероятностью*. Как уже известно, она составляет 20 %.

Предположим, что, внимательно изучив ранее полученные сообщения и частоту появления в них слова «виагра», мы получили дополнительные сведения. Вероятность того, что слово «виагра» присутствовало в

предыдущих спамовых сообщениях, то есть $P(\text{виагра} | \text{спам})$, называется *правдоподобием*. Вероятность того, что слово «виагра» вообще присутствует в каком-либо сообщении, или $P(\text{виагра})$, называется *маргинальным правдоподобием*.

Применяя к этим сведениям теорему Байеса, мы можем вычислить *апостериорную вероятность*, то есть вероятность того, что данное сообщение является спамом. Если апостериорная вероятность больше 50 %, то, скорее всего, сообщение является спамом, а не полезным сообщением и, возможно, его следует отфильтровать. В следующей формуле показано, каким образом теорема Байеса применяется к сведениям, полученным на основании предыдущих сообщений электронной почты:

$$P(\text{спам} | \text{виагра}) = \frac{P(\text{виагра} | \text{спам}) P(\text{спам})}{P(\text{виагра})}$$

Правдоподобие
Априорная вероятность

Апостериорная вероятность
Маргинальное правдоподобие

Для того чтобы вычислить компоненты теоремы Байеса, полезно построить *таблицу частотности (frequency table)* (рис. 4.5, *сверху*), в которой указано, сколько раз слово «виагра» появлялось в спамовых и обычных сообщениях. Как и в случае с перекрестной таблицей, здесь по вертикали указаны уровни переменной класса (спам или не спам), а по горизонтали — уровни признаков (присутствует ли слово «виагра»). В ячейках таблицы указывается количество экземпляров, в которых присутствует сочетание значения класса и значения признака.

Затем таблицу частотности можно использовать для построения *таблицы правдоподобия* (см. рис. 4.5, *снизу*). В строках таблицы правдоподобия указаны значения условной вероятности для слова «виагра» (есть/нет) для сообщений электронной почты, которые являются или не являются спамом.

	Виагра		
Частота	Да	Нет	Всего
Спам	4	16	20
Не спам	1	79	80
Всего	5	95	100

	Виагра		
Правдоподобие	Да	Нет	Всего
Спам	4/20	16/20	20
Не спам	1/80	79/80	80
Всего	5/100	95/100	100

Рис. 4.5. Таблицы частотности и правдоподобия — основа для вычисления апостериорной вероятности того, что письмо является спамом

Как видно из таблицы правдоподобия, $P(\text{виагра} = \text{да} | \text{спам}) = 4 / 20 = 0,20$. Это показывает: вероятность того, что сообщение, которое является

спамом, содержит слово «виагра», составляет 20 %. Кроме того, поскольку $P(A \cap B) = P(B | A) \times P(A)$, то можно рассчитать $P(\text{спам} \cap \text{виагра})$ как $P(\text{виагра} | \text{спам}) \times P(\text{спам}) = (4 / 20) \times (20 / 100) = 0,04$. Тот же результат можно получить и из таблицы частотности, согласно которой четыре из 100 сообщений были спамом, содержащим слово «виагра». В любом случае это в четыре раза больше, чем предыдущая оценка в 0,01, которую мы рассчитали как $P(A \cap B) = P(A) \times P(B)$, ошибочно допустив, что эти события независимы. Это подтверждает важность теоремы Байеса для оценки суммарной вероятности.

Для того чтобы вычислить апостериорную вероятность $P(\text{спам} | \text{виагра})$, мы просто вычислим $P(\text{виагра} | \text{спам}) \times P(\text{спам}) / P(\text{виагра})$, то есть $(4 / 20) \times (20 / 100) / (5 / 100) = 0,80$. Таким образом, вероятность того, что сообщение, в котором содержится слово «виагра», является спамом, составляет 80 %. Исходя из этого результата, любое сообщение, содержащее данное слово, очевидно, следует отфильтровать.

Именно так работают коммерческие спам-фильтры, однако при расчете таблиц частотности и правдоподобия они учитывают гораздо больше слов. В следующем разделе будет показано, как используется этот метод, если задействованы дополнительные признаки.

Наивный байесовский алгоритм

Наивный байесовский алгоритм описывает простой метод применения теоремы Байеса к задачам классификации. Это не единственный метод машинного обучения, в котором используются байесовские методы, однако он наиболее популярен. Своей популярностью этот метод обязан успехам в распознавании текста — в области, где он некогда был негласным стандартом. В табл. 4.2 перечислены достоинства и недостатки этого алгоритма.

Таблица 4.2

Достоинства	Недостатки
<p>Простой, быстрый и очень эффективный.</p> <p>Хорошо работает в условиях зашумленных и отсутствующих данных.</p> <p>Требуется сравнительно мало обучающих примеров, но хорошо работает и с очень большим количеством примеров.</p> <p>Легко получить предположительную вероятность для прогноза</p>	<p>Зачастую опирается на ошибочное предположение о равной важности и независимости признаков.</p> <p>Неидеален для наборов данных с большим количеством числовых признаков.</p> <p>Вычисленные вероятности менее надежны, чем спрогнозированные классы</p>

Наивный байесовский алгоритм получил такое название за то, что он делает некие «наивные» предположения о данных. В частности, наивный байесовский алгоритм предполагает, что все признаки в наборе данных **одинаково важны и независимы друг от друга**. При работе с реальными данными эти предположения редко оказываются верными.

Например, если попытаться идентифицировать спам, отслеживая сообщения электронной почты, то почти наверняка одни признаки окажутся важнее других. Так, отправитель электронной почты может быть более важным признаком спама, чем текст сообщения. Кроме того, слова в теле сообщения не являются независимыми друг от друга: появление некоторых слов — очень хороший показатель того, что в тексте также могут появляться другие слова. Сообщение со словом «виагра» с большой вероятностью будет содержать слова «рецепт» или «наркотики».

Однако в большинстве случаев, даже если эти предположения неверны, наивный байесовский метод все равно достаточно хорошо работает, даже в тех случаях, когда среди признаков обнаруживаются сильные зависимости. Благодаря универсальности и точности этого алгоритма для многих типов условий, особенно для небольших тренировочных наборов данных, наивный байесовский алгоритм часто обоснованно выбирают в качестве первого способа для классификации в обучающих задачах.



Настоящая причина, по которой наивный байесовский алгоритм так хорошо работает, несмотря на ошибочные предположения, была предметом многочисленных спекуляций. Одно из объяснений состоит в том, что не так важно получать точную оценку вероятности, если точны прогнозы. Например, если спам-фильтр правильно идентифицирует спам, так ли уж важно, был ли он уверен в своем прогнозе на 51 % или на 99 %? Одна из дискуссий на эту тему представлена в статье: Domingos P., Pazzani M. On the Optimality of the Simple Bayesian Classifier under Zero-One Loss. *Machine Learning*, 1997. Vol. 29. P. 103–130.

Классификация по наивному байесовскому алгоритму

Расширим наш спам-фильтр, добавив в него еще несколько слов для мониторинга, в дополнение к слову «виагра»: «деньги», «бакалея» и «отписаться». Наивный байесовский алгоритм обучается, составляя таблицу правдоподобия для появления этих четырех слов (обозначены w_1 , w_2 , w_3 и w_4), как показано на рис. 4.6 для 100 электронных писем.

Правдоподобие	Виагра (w_1)		Деньги (w_2)		Бакалея (w_3)		Отписаться (w_4)		Всего
	Да	Нет	Да	Нет	Да	Нет	Да	Нет	
Спам	4/20	16/20	10/20	10/20	0/20	20/20	12/20	8/20	20
Не спам	1/80	79/80	14/80	66/80	8/80	71/80	23/80	57/80	80
Всего	5/100	95/100	24/100	76/100	8/100	91/100	35/100	65/100	100

Рис. 4.6. В расширенной таблице добавлены значения правдоподобия для дополнительных слов в спамовых и обычных сообщениях

По мере поступления новых сообщений необходимо вычислять апостериорную вероятность, чтобы определить, являются ли эти письма

спамом или же нет, с учетом правдоподобия при обнаружении в тексте сообщения соответствующих слов. Предположим, что в сообщении встречаются слова «виагра» и «отписаться», но нет слов «деньги» и «бакалея».

Используя теорему Байеса, мы можем описать задачу, как показано в следующей формуле. Эта формула позволяет вычислить вероятность того, что сообщение является спамом, учитывая, что в нем присутствуют слова «виагра» и «отписаться», а слов «деньги» и «бакалея» там нет:

$$P(\text{спам} | W_1 \cap \neg W_2 \cap \neg W_3 \cap W_4) = \frac{P(W_1 \cap \neg W_2 \cap \neg W_3 \cap W_4 | \text{спам}) P(\text{спам})}{P(W_1 \cap \neg W_2 \cap \neg W_3 \cap W_4)}$$

По ряду причин эту формулу трудно использовать в вычислениях. По мере добавления новых признаков она требует большего объема памяти для хранения вероятностей всех возможных пересекающихся событий. Представьте себе сложность диаграммы Венна для событий, состоящих для четырех слов, не говоря уже о сотнях и более. Многие из потенциальных пересечений никогда не будут обнаружены в прошлых данных, что приведет к нулевой суммарной вероятности и проблемам, которые станут понятны позже.

Чтобы объемы вычислений стали более разумными, можно руководствоваться тем, что наивный байесовский алгоритм делает предположение об отсутствии зависимости между событиями. В частности, он предполагает *условную независимость класса* — это означает, что события являются независимыми, если они обусловлены одним и тем же значением класса. Условная независимость позволяет использовать правило вероятности для независимых событий, согласно которому $P(A \cap B) = P(A) \times P(B)$. Таким образом мы упрощаем числитель, что позволяет умножать отдельные условные вероятности, вместо того чтобы вычислять сложную условную суммарную вероятность.

Наконец, поскольку знаменатель не зависит от целевого класса (спам или не спам), он является постоянным значением и на данный момент может быть проигнорирован. Следовательно, условная вероятность того, что письмо является спамом, может быть описана таким образом:

$$P(\text{спам} | W_1 \cap \neg W_2 \cap \neg W_3 \cap W_4) \propto P(W_1 | \text{спам}) P(\neg W_2 | \text{спам}) P(\neg W_3 | \text{спам}) P(W_4 | \text{спам}) P(\text{спам}).$$

А вероятность того, что сообщение не является спамом, описывается следующей формулой:

$$P(\text{не спам} | W_1 \cap \neg W_2 \cap \neg W_3 \cap W_4) \propto P(W_1 | \text{не спам}) P(\neg W_2 | \text{не спам}) P(\neg W_3 | \text{не спам}) P(W_4 | \text{не спам}) P(\text{не спам}).$$

Обратите внимание, что символ равенства здесь заменен символом пропорциональности (похож на повернутую набок восьмерку с открытым

концом), чтобы подчеркнуть тот факт, что в этой формуле опущен знаменатель.

Используя значения из таблицы правдоподобия, можно подставить числа в эти уравнения. Общее правдоподобие спама составляет:

$$(4 / 20) \times (10 / 20) \times (20 / 20) \times (12 / 20) \times (20 / 100) = 0,012.$$

В то время как вероятность, что письмо не является спамом, равна:

$$(1 / 80) \times (66 / 80) \times (71 / 80) \times (23 / 80) \times (80 / 100) = 0,002.$$

Поскольку $0,012 / 0,002 = 6$, то можно сказать, что сообщение в шесть раз чаще является спамом, чем обычным письмом. Однако, чтобы преобразовать эти числа в вероятности, нужно сделать последний шаг и вновь ввести исключенный знаменатель. По сути, мы должны пересчитать правдоподобие каждого результата, разделив его на общее правдоподобие всех возможных результатов.

Таким образом, вероятность спама равна правдоподобию того, что сообщение является спамом, разделенным на правдоподобие того, что сообщение либо является спамом, либо нет:

$$0,012 / (0,012 + 0,002) = 0,857.$$

Точно так же вероятность того, что сообщение не является спамом, равна правдоподобию того, что сообщение не является спамом, разделенным на вероятность того, что сообщение либо является спамом, либо нет:

$$0,002 / (0,012 + 0,002) = 0,143.$$

Учитывая слова, найденные в этом сообщении, мы ожидаем, что сообщение является спамом с вероятностью 85,7 % и не является спамом с вероятностью 14,3 %. Поскольку это взаимоисключающие и исчерпывающие события, то сумма этих вероятностей равна единице.

Наивный байесовский классификатор, использованный в рассмотренном примере, можно свести к следующей формуле. Вероятность уровня L для класса C , учитывая обоснования, представленные признаками с F_1 по F_n , равна произведению условных вероятностей каждой части обоснования для уровня класса, умноженной на априорную вероятность уровня класса и коэффициент масштабирования $1 / Z$, который преобразует значения правдоподобия в вероятности. Эта формула выглядит следующим образом:

$$P(C_L | F_1 \dots F_n) = \frac{1}{Z} p(C_L) \prod_{i=1}^n p(F_i | C_L).$$

На первый взгляд это уравнение выглядит пугающим, однако, как видно из примера фильтрации спама, последовательность операций довольно проста. Начать следует с построения таблицы частотности, используйте ее для построения таблицы правдоподобия и умножьте условные правдоподобия на «наивное» предположение о независимости признаков. Наконец, разделите результат на суммарное правдоподобие, чтобы преобразовать правдоподобие для каждого класса в вероятность. Когда вы

несколько раз проделаете эти операции вручную, потом все будет получаться на автомате.

Критерий Лапласа

Прежде чем использовать наивный байесовский алгоритм для решения более сложных задач, необходимо учесть некоторые нюансы.

Предположим, мы получили очередное сообщение, в котором на этот раз есть все четыре слова: «виагра», «бакалея», «деньги» и «отписаться». Если мы, как раньше, воспользуемся наивным байесовским алгоритмом и вычислим правдоподобие того, что это сообщение является спамом, то получим:

$$(4 / 20) \times (10 / 20) \times (0 / 20) \times (12 / 20) \times (20 / 100) = 0.$$

А правдоподобие того, что это не спам:

$$(1 / 80) \times (14 / 80) \times (8 / 80) \times (23 / 80) \times (80 / 100) = 0,00005.$$

Следовательно, вероятность спама составляет:

$$0 / (0 + 0,00005) = 0.$$

А вероятность того, что это не спам:

$$0,00005 / (0 + 0,00005) = 1.$$

Согласно этим результатам данное сообщение является спамом с вероятностью 0 % и обычным сообщением — с вероятностью 100 %. Имеет ли смысл такой прогноз? Очевидно, нет. Сообщение содержит несколько слов, обычно связанных со спамом, включая слово «виагра», которое редко используется в обычных сообщениях. Поэтому весьма вероятно, что данное сообщение классифицировано неправильно.

Такие проблемы возникают в тех случаях, когда событие никогда не происходит для одного или нескольких уровней класса, и поэтому совместная вероятность этих событий равна нулю. Например, слово «бакалея» никогда ранее не появлялось в спамовых сообщениях. Следовательно, $P(\text{спам} | \text{бакалея}) = 0 \%$.

Затем, поскольку вероятности в наивной байесовской формуле перемножаются, это нулевое процентное значение приводит к тому, что апостериорная вероятность спама становится равной нулю, давая слову «бакалея» возможность гарантированно аннулировать и отвергать все другие доказательства. Даже если бы без учета этого слова электронное письмо было в подавляющем большинстве случаев расценено как спам, отсутствие слова «бакалея» в спаме всегда наложит вето на все другие доказательства и приведет к тому, что вероятность спама будет равна нулю.

Решение этой проблемы заключается в использовании так называемого *критерия Лапласа*, названного так в честь французского математика Пьера-Симона Лапласа. Критерий Лапласа добавляет к каждому значению в таблице частотности небольшое число, которое

гарантирует, что для каждого класса существует ненулевая вероятность появления каждого признака. Обычно критерию Лапласа присваивается значение, равное 1, — это гарантирует, что каждое сочетание класса и признака будет найдено среди данных хотя бы один раз.



Критерию Лапласа может быть присвоено любое значение, и оно не обязательно должно быть одинаковым для каждого признака. Если мы задались целью дальше совершенствовать байесовские алгоритмы, то могли бы использовать критерий Лапласа, чтобы отразить предполагаемую априорную вероятность того, как каждое свойство относится к каждому классу. На практике, учитывая достаточно большой тренировочный набор данных, это излишне. Поэтому почти всегда используется значение, равное 1.

Посмотрим, как критерий Лапласа повлияет на прогноз для этого сообщения. Используя критерий Лапласа, равный 1, мы прибавим 1 к каждому числителю в функции правдоподобия. Затем нам нужно прибавить 4 к каждому знаменателю условной вероятности в знаменателе, чтобы компенсировать четыре добавленных значения в числителе.

Следовательно, правдоподобие спама равно:

$$(5 / 24) \times (11 / 24) \times (1 / 24) \times (13 / 24) \times (20 / 100) = 0,0004.$$

А правдоподобие того, что сообщение не является спамом:

$$(2 / 84) \times (15 / 84) \times (9 / 84) \times (24 / 84) \times (80 / 100) = 0,0001.$$

Вычисляя $0,0004 / (0,0004 + 0,0001)$, получим, что вероятность спама составляет 80 %, и, следовательно, вероятность того, что это обычное сообщение, составляет около 20 %. Это более правдоподобный результат, чем $P(\text{спам}) = 0$, полученный, когда определяющим фактором послужило одно только слово «бакалея».



Критерий Лапласа был прибавлен к числителю и знаменателю правдоподобий, однако он не был прибавлен к априорным вероятностям — значениям 20/100 и 80/100. Это сделано потому, что наша лучшая оценка суммарной вероятности спама и не спама остается равной 20 и 80 % с учетом того, что наблюдалось в данных.

Использование числовых признаков в наивном байесовском алгоритме

В наивном байесовском алгоритме для изучения данных используются таблицы частотности. Это означает, что каждый признак должен быть категориальным, чтобы создавать комбинации значений классов и признаков, образующих матрицу. Поскольку числовые признаки не имеют

категорий значений, рассмотренный алгоритм не работает с числовыми данными напрямую. Однако существуют способы решить эту проблему.

Одним из простых и эффективных решений является *дискретизация* числовых признаков. Это означает, что числа просто разбиваются на категории, именуемые *интервалами* (bins). Процесс такого разбиения называется *дискретизацией* или *биннингом*. Такой метод лучше всего работает на большом количестве тренировочных данных.

Существует несколько способов дискретизации числовых признаков. Пожалуй, самым распространенным из них является разделение данных на естественные категории, или *точки разделения*, в распределении данных. Например, предположим, что мы ввели в набор данных о спаме признак, в который записывается время суток, когда было отправлено электронное письмо, от 0 до 24 часов. На рис. 4.7 показано, как могут выглядеть сведения о времени, представленные в виде гистограммы.

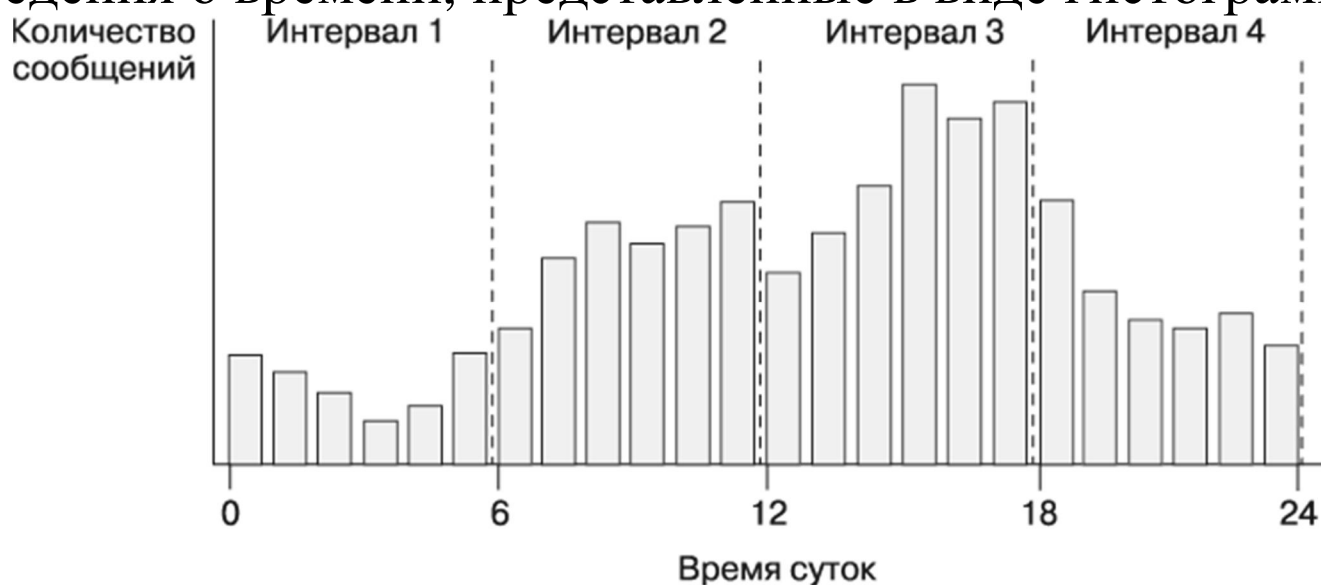


Рис. 4.7. Гистограмма, отображающая распределение времени получения электронных писем

Ранним утром частота сообщений низкая. В рабочее время интенсивность нарастает, а к вечеру снова падает. Так, образуются четыре естественных интервала интенсивности, разделенных пунктирными линиями. Эти линии указывают на точки, в которых числовые данные можно разделить на уровни, чтобы создать новый категориальный признак, который затем можно будет использовать в наивном байесовском алгоритме.

Выбор четырех интервалов был обусловлен естественным распределением данных и догадками о том, как процент спама может изменяться в течение дня. Можно ожидать, что спамеры действуют поздно ночью или же днем, когда люди проверяют электронную почту. Тем не менее, чтобы отразить эти тенденции, мы могли бы так же легко использовать 3 или 12 интервалов.



Если нет явно выраженных точек разделения, одним из вариантов является дискретизация признака с использованием квантилей. Можно разделить

данные на три интервала с помощью тертилей, четыре — с помощью квартилей или пять — с помощью квинтилей.

Следует иметь в виду, что дискретизация числового признака всегда приводит к сокращению информации, поскольку исходная глубина детализации признака сокращается до меньшего числа категорий. Важно соблюсти баланс. Слишком малое количество интервалов может привести к потере важных тенденций. Большое количество интервалов может привести к малым значениям в таблице частотности наивного байесовского алгоритма, что повысит чувствительность алгоритма к зашумленным данным.

Пример: фильтрация спама в мобильном телефоне с помощью наивного байесовского алгоритма

По мере роста популярности мобильных телефонов во всем мире появились новые возможности для распространения рекламы по почте, используемые недобросовестными маркетологами. Такие рекламодатели используют короткие текстовые сообщения (СМС), чтобы привлечь потенциальных потребителей нежелательной рекламой, известной как СМС-спам. Этот тип спама является особенно опасным, поскольку, в отличие от почтового спама, СМС может причинить больше ущерба из-за широкого использования мобильных телефонов. Разработка алгоритма классификации, который бы фильтровал СМС-спам, стала бы полезным инструментом для операторов сотовой связи.

Поскольку наивный байесовский алгоритм успешно применялся для фильтрации спама в электронной почте, вполне вероятно, что он также может быть применен к СМС-спаму. Однако в отличие от спама в электронной почте СМС-спам создает дополнительные проблемы для автоматических фильтров. Размер СМС часто ограничен 160 символами, что сокращает объем текста, по которому можно определить, является ли сообщение нежелательным. Такое ограничение, в сочетании с маленьким размером клавиатуры мобильного телефона, привело к тому, что сформировался своеобразный сокращенный СМС-язык, что еще больше стирает грань между обычными сообщениями и спамом. Посмотрим, каким образом простой наивный байесовский классификатор решает эти проблемы.

Шаг 1. Сбор данных

Для разработки наивного байесовского классификатора воспользуемся данными из набора *SMS Spam Collection*, доступного по адресу <http://www.dt.fee.unicamp.br/~tiago/smsspamcollection/>.



Подробнее о том, как была создана коллекция СМС-спама SMS Spam Collection, читайте в статье: Gómez J.M., Almeida T.A., Yamakami A. On the Validity of a New SMS Spam Collection // Proceedings of the 11th IEEE International Conference on Machine Learning and Applications, 2012.

Этот набор данных включает в себя тексты СМС, а также метки, указывающие на то, является ли данное сообщение нежелательным.

Нежелательные сообщения помечаются как **spam**, а обычные — как **ham**.

В табл. 4.3 показаны некоторые примеры спама и обычных сообщений.

Таблица 4.3

Образцы обычных СМС	Образцы СМС-спама
<p>Лучше. Дождлся пятницы и вчера обожрался, как свинья. Теперь чувствую себя мерзко. Но хотя бы уже ничего не болит.</p> <p>Если он начнет поиск, то получит работу через несколько дней. У него большой потенциал и талант.</p> <p>Я получил другую работу! В больнице, нужно будет анализировать данные или что-то вроде этого, приступаю в понедельник! Понятия не имею, когда закончу диссертацию</p>	<p>Поздравляем! Вы получили 500 ваучеров на CD или 125 подарков и бесплатный инет 2100 нд фото текст МУЗЫКА до 87066.</p> <p>Только декабрь! У тебя был мобильный 11мес+? Тебе предоставлено обновление до новейшей цветной камеры мобильного бесплатно! Для того чтобы получить мобильное обновление, позвони бесплатно на номер 08002986906.</p> <p>Специальное предложение на День святого Валентина! Выиграй более 1000 фунтов стерлингов в нашей викторине и отправься с другом в путешествие на всю жизнь! Отправь GO на 83600 сейчас. 150ф/сообщ получ.</p>

Заметили ли вы какие-либо характерные особенности спама, глядя на эти сообщения? Одна примечательная особенность заключается в том, что в двух из трех спамовых сообщений используется слово «бесплатно», однако это слово не встречается ни в одном из обычных сообщений. И наоборот, в двух обычных сообщениях указаны конкретные дни недели, в то время как ни в одном из спамовых сообщений их нет.

Наивный байесовский классификатор воспользуется преимуществом таких паттернов по частоте слов, чтобы определить, соответствуют ли СМС профилю спама или обычных сообщений. Впрочем, не исключено, что слово «бесплатно» будет появляться и в обычных сообщениях. Вполне вероятно, что в обычных сообщениях будут содержаться дополнительные слова, образующие контекст. Например, в обычном сообщении могут спросить: «Свободны ли вы в воскресенье?» — тогда как в спамовом сообщении может использоваться фраза «бесплатные мелодии».

Классификатор вычислит вероятности спама и обычных сообщений с учетом обстоятельств, представленных всеми словами в сообщении.

Шаг 2. Исследование и подготовка данных

Первый шаг к построению классификатора включает в себя подготовку необработанных данных для анализа. Подготовка текстовых данных — сложная задача, поскольку необходимо преобразовать слова и предложения в форму, понятную компьютеру. Мы преобразуем данные в представление, известное как *неупорядоченный набор слов* (bag-of-words), которое игнорирует последовательность слов и просто предоставляет переменную, указывающую, присутствует ли данное слово в наборе.



Используемый в примере набор данных немного изменен по сравнению с оригиналом, чтобы упростить работу в среде R. Если вы намерены воспроизвести этот пример, скачайте по адресу <https://github.com/dataspelunking/MLwR/> файл `sms_spam.csv` и сохраните его в своем рабочем каталоге R.

Мы начнем с того, что импортируем CSV-данные и сохраним их в виде фрейма данных:

```
> sms_raw <- read.csv("sms_spam.csv",  
stringsAsFactors = FALSE)
```

Воспользовавшись функцией `str()`, мы увидим, что фрейм данных `sms_raw` включает в себя 5559 СМС с двумя признаками: `type` и `text`. Тип СМС закодирован как `ham` или `spam`. В элементе `text` хранится полный текст СМС.

```
> str(sms_raw) 'data.frame':    5559 obs.  
Of    2 variables: $ type: chr  "ham" "ham" "ham"  
"spam" ... $ text: chr  "Hope you are having a  
good week. Just checking in" "K..give back my  
thanks." "Am also doing in cbe only. But have to  
pay." "complimentary 4 STAR Ibiza Holiday or  
£10,000 cash needs your URGENT collection.  
09066364349 NOW from Landline not to lose  
out" |__truncated__ ...
```

Элемент `type` сейчас представляет собой символьный вектор. Поскольку это категориальная переменная, лучше преобразовать ее в коэффициент, как показано в следующем коде:

```
> sms_raw$type <- factor(sms_raw$type)
```

Исследовав результат с помощью функций `str()` и `table()`, мы видим, что теперь `type` соответствующим образом перекодирован как фактор. Кроме того, мы видим, что 747 (около 13 %) СМС из наших данных помечены как `spam`, в то время как остальные помечены как `ham`:

```
> str(sms_raw$type) Factor w/ 2 levels  
"ham", "spam": 1 1 1 2 2 1 1 1 2 1 ...>  
table(sms_raw$type) ham spam4812 747
```

Пока мы оставим текст сообщения в покое. Как вы узнаете из следующего раздела, обработка необработанных СМС потребует использования нового набора мощных инструментов, специально предназначенных для обработки текстовых данных.

Подготовка данных: очистка и стандартизация текстовых данных СМС — это текстовые строки, состоящие из слов, пробелов, чисел и знаков препинания. Обработка сложных данных такого типа требует тщательного продумывания и больших усилий. Необходимо придумать, как убрать цифры и знаки препинания; как обработать неинтересные слова, такие как «и», «но», «или»; как разбить предложения на отдельные слова. К счастью, такая функциональность уже разработана членами R-сообщества и включена в пакет анализа текстовых данных под названием `tm`.



Пакет `tm` был создан Инго Фейнерером (Ingo Feinerer) в процессе написания диссертации в Венском университете экономики и бизнеса. Подробнее об этом пакете читайте статью: Feinerer I., Hornik K., Meyer D. Text Mining Infrastructure in R // Journal of Statistical Software, 2008. Vol. 25. P. 1–54.

Пакет `tm` устанавливается командой `install.packages("tm")` и загружается командой `library(tm)`. Но даже если он уже установлен, возможно, стоит повторить команду установки, чтобы убедиться, что ваша версия обновлена, поскольку пакет `tm` находится в активной разработке. Иногда это приводит к изменениям в его функциональности.



Эта глава была протестирована с использованием версии `tm` 0.7-6, которая была актуальной в феврале 2019 года. Если вы заметите различия в выходных данных или если код не будет работать, возможно, у вас установлена другая версия. В случае если будут замечены значительные изменения, на странице поддержки Racket для этой книги, а также в репозитории GitHub будут опубликованы решения для будущих пакетов `tm`.

Первым шагом при обработке текстовых данных является создание *корпуса* — набора текстовых документов. Документы могут быть короткими или длинными; это может быть набор новостных статей, страниц книги, веб-страниц или даже целых книг. В нашем случае корпус будет представлять собой набор СМС.

Для того чтобы создать корпус, мы воспользуемся функцией `VCorpus()` из пакета `tm`, название которой означает *volatile corpus* — «временный корпус» — временный, поскольку хранится в памяти, а не на диске (для доступа к постоянному корпусу, который хранится в базе данных, используется функция `PCorpus()`).

Функции `VCorpus()` необходимо указать источник документов для

создания корпуса: это может быть файловая система компьютера, база данных, сеть или что-либо еще.

Поскольку мы уже загрузили в среду R текст СМС, то для создания исходного объекта из существующего вектора `sms_raw$text` мы будем использовать функцию чтения `VectorSource()`, а результат затем передадим в `VCorpus()` следующим образом:

```
> sms_corpus <-  
VCorpus(VectorSource(sms_raw$text))
```

Полученный объект корпуса сохраняется под именем `sms_corpus`.



Указав дополнительный параметр `readerControl`, можно использовать функцию `VCorpus()` для импорта текста из таких источников, как PDF-файлы и файлы Microsoft Word. Подробнее об этом читайте в разделе `Data Import` в кратком описании пакета `tm`, которое доступно с помощью команды `vignette("tm")`.

При выводе корпуса на печать мы видим, что он содержит документы для каждого из 5559 СМС, входящих в тренировочный набор данных:

```
>  
print(sms_corpus) <<VCorpus>>Metadata: corpus  
specific: 0, document level (indexed):  
0Content: documents: 5559
```

Теперь, поскольку корпус, составленный с помощью `tm`, представляет собой, в сущности, сложный список, для выбора документов в корпусе мы можем использовать операции для работы со списками. Для просмотра сводного результата можно воспользоваться функцией `inspect()`. Например, следующая команда позволяет увидеть сводную информацию для первого и второго СМС в корпусе:

```
>  
inspect(sms_corpus[1:2]) <<VCorpus>>Metadata: co  
rpus specific: 0, document level (indexed):  
0Content: documents:  
2[[1]] <<PlainTextDocument>>Metadata: 7Content:  
chars:  
49[[2]] <<PlainTextDocument>>Metadata: 7Content:  
chars: 23
```

Для того чтобы просмотреть реальный текст сообщения, нужно применить к этим сообщениям функцию `as.character()`. Чтобы просмотреть отдельное сообщение, используйте

функцию `as.character()` для одного элемента списка и не забудьте о двойных скобках:

```
> as.character(sms_corpus[[1]])[1] "Hope you are having a good week. Just checking in"
```

Чтобы просмотреть несколько документов, нужно применить `as.character()` к нескольким элементам в объекте `sms_corpus`. Воспользуемся функцией `lapply()` — одной из семейства R-функций, которые выполняют действия с каждым элементом структуры данных R. Это семейство функций, в которое, кроме прочих, также входят функции `apply()` и `sapply()`, относящиеся к ключевым идиомам языка R. Опытные R-кодеры используют их так же, как циклы `for` и `while` в других языках программирования, поскольку эти функции позволяют писать более читабельный (а иногда и более эффективный) код. Функция `lapply()`, примененная вместе с `as.character()` к подмножеству элементов корпуса, выглядит следующим образом:

```
> lapply(sms_corpus[1:2], as.character)$'1'[1]
"Hope you are having a good week. Just checking in"
$'2'[1] "К..give back my thanks."
```

Как уже отмечалось, корпус содержит необработанный текст 5559 сообщений. Для того чтобы выполнить анализ, нам нужно разбить эти сообщения на отдельные слова, но предварительно следует очистить текст, чтобы стандартизировать слова и удалить знаки препинания, которые зашумляют результат. Например, нужно, чтобы строки *Hello!*, *HELLO* и *hello* считались вариантами одного и того же слова.

Функция `tm_map()` предоставляет метод преобразования корпуса `tm` (также известного как отображение). Мы будем использовать эту функцию, чтобы очистить корпус, применяя серию преобразований, и сохраним результат в новом объекте с именем `corpus_clean`.

Первым преобразованием будет стандартизация сообщений таким образом, чтобы в них присутствовали только строчные символы. Для этого в R есть функция `tolower()`, которая возвращает текстовую строку, все буквы которой преобразованы в строчные. Для того чтобы применить эту функцию к корпусу, нужно использовать функцию-оболочку `content_transformer()` из пакета `tm`, чтобы функция `tolower()` рассматривалась как функция преобразования, применяемая к корпусу. Полностью команда выглядит следующим образом:

```
> sms_corpus_clean <-  
tm_map(sms_corpus, content_transformer(tolower  
) )
```

Чтобы убедиться, что команда сработала должным образом, возьмем первое сообщение в исходном корпусе и сравним его с тем же сообщением в преобразованном корпусе:

```
> as.character(sms_corpus[[1]])[1] "Hope you  
are having a good week. Just checking in">  
as.character(sms_corpus_clean[[1]])[1] "hope you  
are having a good week. just checking in"
```

Как и предполагалось, прописные буквы в обработанном корпусе заменены строчными.



Для более сложных процессов обработки и очистки текста, таких как паттерн сопоставления и замены `grep`, может использоваться функция `content_transformer()`. Для этого просто напишите свою функцию и перед применением оберните ее в функцию `tm_map()`.

Продолжим очистку данных, удалив из СМС числа. Некоторые числа могут содержать полезную информацию, однако большинство из них, скорее всего, будут уникальными для конкретных отправителей и, следовательно, не предоставят полезные паттерны для всех сообщений.

С учетом этого мы удалим из корпуса все числа следующим образом:

```
> sms_corpus_clean <- tm_map(sms_corpus_clean,  
removeNumbers )
```



Обратите внимание, что в предыдущем коде не использовалась функция `content_transformer()`. Дело в том, что функция `removeNumbers()` встроена в `tm`, как и некоторые другие функции отображения, которые не нужно обертывать. Чтобы увидеть, какие еще преобразования встроены в пакет `tm`, просто введите команду `getTransformations()`.

Следующая задача — удалить из СМС такие слова-связки, как «к», «и», «но», «или», которые называются *стоп-словами*. Обычно они удаляются перед анализом текстовых данных. Это связано с тем, что, хотя эти слова встречаются очень часто, они не дают полезной для машинного обучения информации.

Вместо того чтобы самостоятельно определять список стоп-слов, мы воспользуемся функцией `stopwords()`, предоставляемой пакетом `tm`. Эта функция позволяет получить доступ к наборам стоп-слов на разных языках. По умолчанию используются стоп-слова на английском. Для

просмотра списка стоп-слов по умолчанию введите в командной строке R команду `stopwords()`. Для просмотра других доступных языков и параметров введите `?stopwords`, чтобы открыть страницу документации.



Даже в пределах одного языка не существует единого универсального списка стоп-слов. Например, список стоп-слов английского языка, включенный в `tm` по умолчанию, насчитывает около 174 слов, а другой его вариант включает в себя 571 слово. Вы даже можете указать свой собственный список стоп-слов. Независимо от того, какой список вы предпочтете, помните о цели этого преобразования, которая заключается в том, чтобы исключить лишние данные, но сохранить как можно больше полезной информации.

Нам нужен способ, с помощью которого можно удалить все слова, появляющиеся в списке стоп-слов. Решение заключается в функции `removeWords()`, которая является преобразованием, включенным в пакет `tm`. Как и ранее, мы воспользуемся функцией `tm_map()`, чтобы применить это отображение к данным, подставляя функцию `stopwords()` в качестве параметра, чтобы точно указать, какие именно слова мы бы хотели удалить. Полная команда выглядит следующим образом:

```
> sms_corpus_clean <-  
tm_map(sms_corpus_clean, removeWords,  
stopwords())
```

Поскольку функция `stopwords()` просто возвращает вектор стоп-слов, то при желании мы могли бы заменить вызов этой функции на собственный вектор слов, подлежащий удалению. Таким образом можно было бы расширить либо уменьшить список стоп-слов по своему желанию или же удалить совсем другой набор слов.

Продолжая процесс очистки, мы также можем исключить из текстовых сообщений знаки пунктуации с помощью встроенного преобразования `removePunctuation()`:

```
> sms_corpus_clean <- tm_map(sms_corpus_clean,  
removePunctuation)
```

Преобразование `removePunctuation()` полностью удаляет из текста знаки препинания, что может привести к непредвиденным последствиям. Например, рассмотрим, что произойдет, если это преобразование применяется следующим образом:

```
> removePunctuation("hello...world")[1]  
"helloworld"
```

Как видите, отсутствие пробела после скобок привело к тому, что слова *hello* и *world* слились в одно. Сейчас это не станет существенной проблемой, однако в будущем следует иметь такое в виду.



Для того чтобы обойти это стандартное поведение `removePunctuation()`, можно создать свою функцию, которая не удаляет знаки пунктуации, а заменяет их:

```
> replacePunctuation <- function(x) {  
  gsub("[:,punct:"]+", " ", x)  
}
```

Здесь R-функция `gsub()` используется для замены любых знаков пунктуации в тексте `x` пробелом. Затем можно задействовать функцию `replacePunctuation()` совместно с `tm_map()`, как и в случае с другими преобразованиями.

Еще одна распространенная стандартизация для текстовых данных заключается в приведении слов к их корневой форме — этот процесс называется *морфологическим поиском* (*stemming*). В процессе морфологического поиска берутся такие слова, как *learned*, *learning* или *learns*, и путем удаления суффикса преобразуются в базовую форму *learn*. Это позволяет алгоритмам машинного обучения рассматривать связанные слова как единую концепцию, вместо того чтобы пытаться создать паттерн для каждого варианта.

В пакете `tm` морфологический поиск реализован путем интеграции с пакетом `SnowballC`. На момент написания этой книги `SnowballC` по умолчанию не устанавливался с `tm`, поэтому, если вы еще его не установили, сделайте это с помощью команды `install.packages("SnowballC")`.



Пакет `SnowballC` поддерживается Миланом Буше-Валатом (Milan Bouchet-Valat) и обеспечивает R-интерфейс для библиотеки `libstemmer` на основе языка C, которая основана на разработанном М.Ф. Портером (M. F. Porter) алгоритме морфологического поиска `Snowball` — широко используемом методе морфологического поиска с открытым исходным кодом. Подробнее об этом алгоритме читайте на сайте <http://snowballstem.org>.

В состав пакета `SnowballC` входит функция `wordStem()`, которая принимает символьный вектор и возвращает тот же вектор слов в их начальной форме. Например, эта функция правильно анализирует предложенные ранее варианты слова *learn*:


```
> library(SnowballC) > wordStem(c("learn",
"learned", "learning", "learns"))[1]
"learn" "learn" "learn" "learn"
```

Для того чтобы применить функцию `wordStem()` ко всему корпусу текстовых документов, в пакете `tm` есть преобразование `stemDocument()`. Мы применим это преобразование к корпусу с помощью функции `tm_map()` точно так же, как и раньше:

```
> sms_corpus_clean <- tm_map(sms_corpus_clean,
stemDocument)
```



Если при применении преобразования `stemDocument()` вы получите сообщение об ошибке, убедитесь, что у вас установлен пакет `SnowballC`. Если пакет установлен, но при этом появится сообщение `All scheduled cores encountered errors` (Ошибка во всех запланированных ядрах), можете попытаться принудительно выполнить команду `tm_map()` к одному ядру, указав дополнительный параметр `mc.cores = 1`.

После удаления чисел, стоп-слов и знаков препинания, а также выполнения морфологического поиска мы получим текстовые сообщения с пробелами, которые когда-то разделяли удаленные фрагменты. Поэтому последним шагом в процессе очистки текста будет удаление лишних пробелов с помощью встроенного преобразования `stripWhitespace()`:

```
> sms_corpus_clean <- tm_map(sms_corpus_clean,
stripWhitespace)
```

В следующей таблице показаны первые три сообщения в корпусе СМС до и после очистки. В сообщениях остались наиболее интересные слова, а пунктуация и заглавные буквы удалены (табл. 4.4).

Таблица 4.4

СМС до очистки	СМС после очистки
<pre>> as.character(sms_corpus[1:3]) [[1]] Hope you are having a good week. Just checking in [[2]] K..give back my thanks. [[3]] Am also doing in cbe only. But have to pay.</pre>	<pre>> as.character(sms_corpus_clean[1:3]) [[1]] hope good week just check [[2]] kgive back thank [[3]] also cbe pay</pre>

Подготовка данных: разбиение текстовых документов на отдельные слова

Теперь, когда данные обработаны так, как нам надо, остался последний шаг — разделить сообщения на отдельные слова путем *токенизации*. Токен —

это отдельный элемент текстовой строки; в данном случае токены являются словами.

Как вы могли догадаться, пакет `tm` предоставляет функциональные возможности для токенизации корпуса СМС.

Функция `DocumentTermMatrix()` принимает на входе корпус и создает структуру данных, называемую *матрицей «документ — термин»* (Document-Term Matrix, DTM), строки которой соответствуют документам (СМС), а столбцы — терминам (словам).



Пакет `tm` также предоставляет структуру данных для матрицы «документ — термин», которая представляет собой обычную транспонированную матрицу DTM. Строки в ней соответствуют терминам, а столбцы — документам. Зачем нужны обе эти матрицы? Иногда удобнее работать с одной, иногда — с другой. Например, если документов мало, а список слов велик, то имеет смысл использовать TDM, потому что обычно проще отображать много строк, чем много столбцов. В любом случае эти две матрицы, как правило, взаимозаменяемы.

В каждой ячейке матрицы хранится число, соответствующее тому, сколько раз слово из столбца встречается в документе (отображается в строке). На рис. 4.8 показана небольшая часть DTM для корпуса СМС: полная матрица содержит 5559 строк и более 7000 столбцов.

Номер сообщения	baloon	balls	bam	bambling	band
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0
5	0	0	0	0	0

Рис. 4.8. DTM для СМС в основном заполнена нулями

Если все ячейки в таблице равны нулю, это значит, что ни одно из слов, перечисленных в заголовках столбцов, не встречается ни в одном из первых пяти сообщений в корпусе. Отсюда и название структуры данных — *разреженная матрица*; подавляющее большинство ячеек этой матрицы заполнено нулями. В реальных условиях это означает: несмотря на то что каждое сообщение должно содержать хотя бы одно слово, вероятность того, что некое конкретное слово появится в данном сообщении, мала.

Для создания разреженной DTM-матрицы из корпуса данных используется одна команда из пакета `tm`:

```
> sms_dtm <-  
DocumentTermMatrix(sms_corpus_clean)
```

Эта команда создаст объект `sms_dtm`, который содержит токенизированный корпус с параметрами по умолчанию, что требует минимальной обработки. Параметры по умолчанию нам подходят, поскольку мы уже подготовили корпус вручную.

Однако если бы мы не выполнили предварительную обработку данных, то мы могли бы сделать это сейчас, предоставив список управляющих параметров `control` для переопределения значений по умолчанию. Например, чтобы создать DTM непосредственно из необработанного, сырого СМС-корпуса, можно использовать следующую команду:

```
> sms_dtm2 <- DocumentTermMatrix(sms_corpus,
control = list(tolower =
TRUE, removeNumbers = TRUE, stopwords =
TRUE, removePunctuation = TRUE, stemming =
TRUE))
```

Эта команда выполняет те же операции предварительной обработки для корпуса СМС и в том же порядке, которые были описаны ранее. Однако, сравнив `sms_dtm` с `sms_dtm2`, мы заметим небольшую разницу в количестве элементов матрицы:

```
> sms_dtm<<DocumentTermMatrix (documents: 5559,
terms: 6559)>>Non-/sparse entries:
42147/36419334Sparsity : 100%Maximal
term length: 40Weighting : term
frequency (tf)> sms_dtm2<<DocumentTermMatrix
(documents: 5559, terms: 6961)>>Non-/sparse
entries: 43221/38652978Sparsity :
100%Maximal term length: 40Weighting :
term frequency (tf)
```

Причина такого несоответствия вызвана незначительным различием в последовательности выполнения этапов предварительной обработки. Функция `DocumentTermMatrix()` выполняет операции очистки к текстовым строкам после их разделения на слова. Таким образом, в ней используется несколько другой метод удаления стоп-слов. Следовательно, некоторые слова разделяются иначе, чем когда они очищаются перед токенизацией.



Для того чтобы две предыдущие DTM-матрицы были идентичными, можно заменить примененную по умолчанию функцию стоп-слов собственной, которая использует оригинальную функцию замены. Нужно просто заменить `stopwords = TRUE` на следующее:

```
stopwords = function(x) { removeWords(x,  
stopwords()) }
```

Различия между этими двумя вариантами иллюстрируют важный принцип очистки текстовых данных: последовательность операций имеет значение. Зная об этом, очень важно продумать, как ранние этапы процесса повлияют на последующие. Представленная здесь последовательность, как правило, будет работать, однако если необходимо более тщательно адаптировать процесс к конкретным наборам данных и случаям использования, то может потребоваться изменение этой последовательности. Например, если есть определенные слова, которые желательно исключить из матрицы, подумайте, следует искать их до или после морфологического поиска. Подумайте также о том, как исключение знаков препинания — удаление или замена пробелами — повлияет на эти процессы.

Подготовка данных: создание тренировочного и тестового набора данных

Теперь, когда данные готовы для анализа, необходимо разделить их на тренировочный и тестовый наборы данных, чтобы созданный классификатор нежелательной почты можно было проверить на тех данных, которые ему ранее не встречались. Но, даже несмотря на то, что мы должны скрыть от классификатора содержимое тестового набора данных, важно, чтобы разделение данных происходило после их очистки и обработки. И тренировочные, и тестовые данные должны пройти одни и те же подготовительные этапы.

Разделим данные на две части: 75 % для обучения и 25 % для тестирования. Поскольку СМС представлены в случайном порядке, можно просто взять первые 4169 сообщений для обучения и оставить 1390 для тестирования. К счастью, DTM-объект очень похож на фрейм данных и может быть разделен с помощью стандартных операций `[row, col]`. Поскольку в нашем DTM-объекте СМС представлены в виде строк, а слова — в виде столбцов, следует запросить определенный диапазон строк и все столбцы для каждого набора данных:

```
> sms_dtm_train <- sms_dtm[1:4169, ]>  
sms_dtm_test <- sms_dtm[4170:5559, ]
```

Для удобства полезно будет сохранить и пару векторов с метками для каждой строки в тренировочной и тестовой матрицах. Эти метки не хранятся в DTM, поэтому их нужно извлечь из исходного фрейма данных `sms_raw`:

```
> sms_train_labels <- sms_raw[1:4169, ]$type>  
sms_test_labels <- sms_raw[4170:5559, ]$type
```

Для того чтобы убедиться, что полученные подмножества являются такими же репрезентативными, что и полный набор данных об СМС, сравним долю спама в тренировочном и тестовом фреймах данных:

```
>
prop.table(table(sms_train_labels))      ham
spam0.8647158      0.1352842>
prop.table(table(sms_test_labels))      ham
spam0.8683453      0.1316547
```

И тренировочные, и тестовые данные содержат около 13 % спама. Это говорит о том, что спамовые сообщения равномерно распределены между обоими наборами данных.

Визуализация текстовых данных: облака слов

Облако слов — это вариант визуального представления частоты появления слов в текстовых данных. Облако состоит из слов, случайно разбросанных по какой-то области. Слова, которые встречаются в тексте чаще, отображаются более крупным шрифтом, а менее распространенные — более мелким. Такой тип представления стал популярен для отслеживания актуальных тем на сайтах социальных сетей.

В пакете `wordcloud` есть простая R-функция для построения диаграмм подобного типа. Воспользуемся ею для визуализации слов в СМС. Сравнение облаков слов для спамовых и обычных сообщений поможет определить, насколько успешным является спам-фильтр, построенный на основе наивного байесовского алгоритма. Если вы еще этого не сделали, установите и загрузите пакет `wordcloud`, введя в командной строке R

команды `install.packages("wordcloud")` и `library(wordcloud)`.



Пакет `wordcloud` написан Яном Феллоузом (Ian Fellows). Подробнее об этом пакете читайте в блоге автора по адресу <http://blog.fellstat.com/?cat=11>. Облако слов можно создать непосредственно из корпуса `tm` с помощью следующего синтаксиса:

```
> wordcloud(sms_corpus_clean, min.freq = 50,
random.order = FALSE)
```

В результате из подготовленного СМС-корпуса будет создано облако слов. Поскольку мы указали параметр `random.order=FALSE`, облако будет располагаться не в случайном порядке — слова с более высокой частотой появления разместятся ближе к центру. Если не

указать `random.order`, то по умолчанию слова в облаке будут располагаться случайным образом.

Параметр `min.freq` указывает на то, сколько раз слово должно появиться в корпусе, чтобы отображаться в облаке. Поскольку частота, равная 50, составляет около 1 % от корпуса, это означает, что для включения в облако слово должно присутствовать как минимум в 1 % СМС.



Вы можете получить предупреждающее сообщение о том, что R не смог уместить все слова на рисунке. Если это произойдет, попробуйте увеличить `min.freq`, чтобы уменьшить количество слов в облаке. Может также помочь использование параметра `scale`, который уменьшает размер шрифта.

Полученное облако слов должно выглядеть примерно так, как показано на рис. 4.9.



Рис. 4.9. Облако слов, отображающее сл

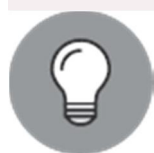
Возможно, мы получили бы более интересную визуализацию, если бы сравнили облака слов для СМС-спама и обычных сообщений. Поскольку мы не создавали отдельные корпуса для спама и не спама, сейчас самое время обратить внимание на очень полезное свойство функции `wordcloud()`. Принимая вектор необработанных текстовых строк, она автоматически применяет к нему общие процессы подготовки текста и затем отображает облако слов.

Воспользуемся R-функцией `subset()` и в результате получим подмножество данных `sms_raw`, отобранное по типу СМС. Сначала мы создадим подмножество, в котором значением `type` является `spam`:

```
> spam <- subset(sms_raw, type == "spam")
```

Затем сделаем то же самое для подмножества `ham`:

```
> ham <- subset(sms_raw, type == "ham")
```



Обратите внимание: здесь стоит двойной знак равенства. Подобно многим языкам программирования, в R знак `==` используется для проверки равенства. Если вы случайно поставите один знак равенства, то получите намного большее подмножество, чем ожидали!

Теперь у нас есть два фрейма данных, `spam` и `ham`, в каждом из которых признак `text` содержит необработанные текстовые строки СМС. Создать из них облака слов так же просто, как и раньше. На сей раз мы воспользуемся параметром `max.words`, чтобы получить 40 самых распространенных слов в каждом из двух наборов.

Параметр `scale` позволяет выбрать максимальный и минимальный размер шрифта для слов в облаке. Всегда настраивайте эти параметры по своему усмотрению, как показано в следующем коде:

```
> wordcloud(spam$text, max.words = 40, scale =  
c(3, 0.5))> wordcloud(ham$text, max.words = 40,  
scale = c(3, 0.5))
```

Полученные облака слов должны выглядеть примерно так, как указано на рис. 4.10.

Вы уже догадались, какое из этих облаков соответствует спаму, а какое — обычным сообщениям?



Рис. 4.10. Расположенные рядом облака слов, соответствующие СМС-спаму и обычным сообщениям



Из-за рандомизации у вас облака слов могут выглядеть немного иначе. Запустив функцию `wordcloud()` несколько раз, вы можете выбрать то облако, которое лучше всего подходит для презентации. Как вы, наверное, догадались, облако спама находится слева. Спамовые сообщения включают в себя такие слова, как *urgent* («срочно»), *free* («бесплатно»), *mobile* («мобильный»), *claim* («требование») и *stop* («остановитесь»); в облаке для обычных сообщений эти слова не появляются вообще. Вместо этого в обычных сообщениях встречаются такие слова, как *can* («можете»), *sorry* («извините»), *need* («нужно») и *time* («время»). Эти резкие различия свидетельствуют о том, что в нашей наивной байесовской модели будут содержаться некоторые ключевые слова, позволяющие различать классы СМС.

Подготовка данных: создание признаков-индикаторов для часто встречающихся слов

Последним шагом в процессе подготовки данных является преобразование разреженной матрицы в структуру данных, которую можно использовать для обучения наивного байесовского классификатора. Сейчас наша разреженная матрица насчитывает более 6500 признаков: по одному для каждого слова, которое встречается хотя бы в одном СМС. Вряд ли все это будет полезно для классификации. Чтобы уменьшить количество признаков, исключим все слова, которые встречаются менее чем в пяти сообщениях или менее чем в 0,1 % записей в тренировочных данных.

Для поиска часто встречающихся слов воспользуемся функцией `findFreqTerms()` из пакета `tm`. Она принимает DTM-объект и возвращает вектор символов, содержащий слова, которые встречаются не реже чем заданное минимальное количество раз.

Например, следующая команда отображает слова, встречающиеся в матрице `sms_dtm_train` как минимум пять раз:

```
> findFreqTerms(sms_dtm_train, 5)
```

Результатом выполнения функции является символьный вектор, поэтому сохраним наши часто встречающиеся слова для дальнейшего использования:

```
> sms_freq_words <-  
findFreqTerms(sms_dtm_train, 5)
```

Просмотрев содержимое вектора, мы увидим, что 1139 слов встречаются как минимум в пяти СМС:

```
> str(sms_freq_words)chr [1:1139] "£wk" "€ ã"  
"€ ã" "abiola" "abl" "abt" "accept" "access"  
"account" "across" "act" "activ" ...
```

Теперь нужно отфильтровать DTM-объект так, чтобы включить в него только те слова, которые присутствуют в векторе часто встречающихся слов. Как и прежде, для получения определенных разделов DTM мы будем использовать операции в стиле фрейма данных `[row, col]`. Не забывайте, что столбцы соответствуют словам, содержащимся в DTM. Этим фактом можно воспользоваться, чтобы ограничить DTM-объект конкретными словами. Поскольку нам нужны все строки и только те столбцы, которые соответствуют словам из вектора `sms_freq_words`, то команды будут выглядеть так:

```
> sms_dtm_freq_train <- sms_dtm_train[ ,  
sms_freq_words]> sms_dtm_freq_test <-  
sms_dtm_test[ , sms_freq_words]
```

Тренировочный и тестовый наборы данных теперь включают в себя 1139 признаков, которые соответствуют словам, фигурирующим как минимум в пяти сообщениях.

Наивный байесовский классификатор обычно обучается на данных с категориальными признаками. Это усложняет дело, поскольку в ячейках разреженной матрицы хранятся числовые значения, которые означают то, сколько раз данное слово встречается в сообщении. Нам нужно заменить их категориальной переменной, которая принимала бы только значения «да» (`Yes`) или «нет» (`No`), в зависимости от того, появляется ли данное слово в сообщении.

Следующая команда определяет функцию `convert_counts()` для преобразования значений в строки `Yes` или `No`:

```
> convert_counts <- function(x) { x <-  
ifelse(x > 0, "Yes", "No") }
```

Некоторые составляющие этой функции уже знакомы вам. Первая строка определяет функцию.

Оператор `ifelse(x>0, "Yes", "No")` преобразует значение `x` так, что если оно больше 0, то будет заменено на `"Yes"`, а если нет — то на `"No"`. Наконец, в последней строке возвращается преобразованный вектор `x`.

Теперь нужно применить `convert_counts()` ко всем столбцам нашей разреженной матрицы. Как можно догадаться, R-функция, позволяющая сделать это, называется `apply()` и используется так же, как ранее использовалась `lapply()`.

Функция `apply()` помогает задействовать функцию для каждой строки или столбца в матрице. Ее параметр `MARGIN` позволяет выбрать строки или столбцы. Мы воспользуемся значением `MARGIN=2`, так как нас интересуют столбцы (`MARGIN=1` используется для строк). Команды для преобразования тренировочной и тестовой матриц выглядят следующим образом:

```
> sms_train <- apply(sms_dtm_freq_train, MARGIN  
= 2, convert_counts)> sms_test <-  
apply(sms_dtm_freq_test, MARGIN =  
2, convert_counts)
```

В результате получим две матрицы символьного типа, в ячейках каждой из которых будут содержаться значения `"Yes"` или `"No"`, в зависимости от того, появляется ли слово, представленное столбцом, в любом месте сообщения, представленного строкой.

Шаг 3. Обучение модели на данных

После того как исходные СМС преобразованы в формат, который может быть представлен в виде статистической модели, пора применить наивный байесовский алгоритм. Для оценки вероятности того, что данное СМС является спамом, алгоритм будет использовать информацию о наличии или отсутствии в сообщении определенных слов.

Реализация наивного байесовского алгоритма, которую мы будем использовать, входит в состав пакета `e1071`. Этот пакет разработан на факультете статистики Венского технологического университета (TU Wien) и включает в себя различные функции для машинного обучения. Прежде чем продолжить работу, установите и загрузите этот пакет (если вы этого еще не сделали) с помощью команд `install.packages("e1071")` и `library(e1071)`.



Многие технологии машинного обучения реализованы в нескольких R-пакетах, и наивный байесовский алгоритм не является исключением. Еще

один его вариант, `NaiveBayes()`, входящий в состав пакета `klaR`, практически идентичен тому, что реализован в пакете `e1071`. Используйте тот вариант, который вам больше нравится.

В отличие от алгоритма `k-NN`, который мы применяли для классификации в предыдущей главе, обучение наивного байесовского алгоритма и его использование для классификации происходят на разных этапах. Тем не менее, как показано далее, эти этапы весьма просты.

Синтаксис наивного байесовского классификатора
Использование функции <code>NaiveBayes()</code> из пакета <code>e1071</code>
Построение классификатора: <pre>m <- naiveBayes(train, class, laplace = 0)</pre> <p><code>train</code> – фрейм данных или матрица с тренировочными данными; <code>class</code> – факторный вектор, содержащий классы для всех строк тренировочных данных; <code>laplace</code> – число, определяющее критерий Лапласа (по умолчанию 0).</p> <p>Функция возвращает объект наивной байесовской модели, который может быть использован для дальнейших прогнозов.</p>
Прогнозирование: <pre>p <- predict(m, test, type = "class")</pre> <p><code>m</code> – модель, обученная с помощью функции <code>naiveBayes()</code>; <code>test</code> – фрейм данных или матрица, содержащая тестовые данные с теми же признаками, что и тренировочные данные, использованные при построении классификатора; <code>type</code> – принимает значение <code>"class"</code> или <code>"row"</code>; определяет, в какой форме должен выдаваться прогноз: в виде значения класса или вероятности спрогнозированной строки.</p> <p>Эта функция возвращает вектор спрогнозированных значений класса или вероятностей, в зависимости от значения параметра <code>type</code>.</p>
Пример: <pre>sms_classifier <- naiveBayes(sms_train, sms_type) sms_predictions <- predict(sms_classifier, sms_test)</pre>

Для того чтобы построить модель на основе матрицы `sms_train`, воспользуемся следующей командой:

```
> sms_classifier <- naiveBayes(sms_train,  
sms_train_labels)
```

Теперь переменная `sms_classifier` содержит объект классификатора `naiveBayes`, который можно использовать для прогнозирования.

Шаг 4. Оценка эффективности модели

Чтобы оценить СМС-классификатор, нам нужно проверить его прогнозы на сообщениях, не входящих в тестовый набор данных. Напомню, что

признаки этих сообщений хранятся в матрице `sms_test`, а метки классов (`spam` или `ham`) — в векторе `sms_test_labels`. Обученному классификатору присвоено имя `sms_classifier`. Используем его для генерации прогнозов, а затем сравним прогнозные значения с истинными.

Для прогнозирования используется функция `predict()`. Будем хранить спрогнозированные значения в векторе `sms_test_pred`. Предоставим этой функции имена нашего классификатора и тестового набора данных, как показано ниже:

```
> sms_test_pred <- predict(sms_classifier,
sms_test)
```

Для того чтобы сравнить прогнозы с истинными значениями, как и раньше, воспользуемся функцией `CrossTable()` из пакета `gmodels`. На этот раз мы введем некоторые дополнительные параметры, чтобы исключить ненужные пропорции ячеек, и используем параметр `dnn` (названия измерений) для переименования строк и столбцов, как показано в следующем коде:

```
> library(gmodels) > CrossTable(sms_test_pred,
sms_test_labels, prop.chisq = FALSE, prop.c =
FALSE, prop.r = FALSE, dnn = c('predicted',
'actual'))
```

В результате получим таблицу:

Total Observations in Table: 1390

predicted	actual		Row Total
	ham	spam	
ham	1201 0.864	30 0.022	1231
spam	6 0.004	153 0.110	159
Column Total	1207	183	1390

Как видно из таблицы, из 1390 СМС всего $6 + 30 = 36$ (2,6 %) были классифицированы ошибочно. Среди ошибок — 6 из 1207 обычных сообщений, которые были идентифицированы как спам, и 30 из 183 спамовых сообщений, которые были неправильно помечены как обычные. Учитывая те незначительные усилия, которые мы потратили на этот проект, такой уровень эффективности весьма впечатляет. Это исследование наглядно показывает, почему наивный байесовский метод так часто используется для классификации текста: он работает на удивление хорошо без какой-либо настройки.

Впрочем, шесть обычных сообщений, которые были ошибочно классифицированы как спам, могут вызвать серьезные проблемы при развертывании алгоритма фильтрации, поскольку фильтр может привести

к тому, что человек пропустит важное текстовое сообщение. Необходимо выяснить, можно ли немного улучшить модель, чтобы повысить ее эффективность.

Шаг 5. Повышение эффективности модели

Возможно, вы заметили, что при обучении модели не были заданы значения критерия Лапласа. Благодаря этому слова, не появившиеся в спаме или обычных сообщениях, тоже имеют неоспоримое право голоса в процессе классификации. То, что в тренировочных данных слово «рингтон» появилось только в спамовых сообщениях, не означает, что любое сообщение с этим словом следует классифицировать как спам.

Построим наивную байесовскую модель, но на этот раз с параметром `laplace=1`:

```
> sms_classifier2 <- naiveBayes(sms_train,
sms_train_labels, laplace = 1)
```

Затем сделаем прогнозы:

```
> sms_test_pred2 <- predict(sms_classifier2,
sms_test)
```

А теперь сравним прогнозируемые классы с реальной классификацией, используя перекрестную таблицу:

```
> CrossTable(sms_test_pred2,
sms_test_labels, prop.chisq = FALSE, prop.c =
FALSE, prop.r = FALSE, dnn = c('predicted',
'actual'))
```

В результате получим таблицу:

Total Observations in Table: 1390

predicted	actual		Row Total
	ham	spam	
ham	1202 0.996	28 0.153	1230
spam	5 0.004	155 0.847	160
Column Total	1207 0.868	183 0.132	1390

Добавление критерия Лапласа уменьшило количество ложных положительных срабатываний (обычных сообщений, ошибочно классифицированных как спам) с шести до пяти, а количество ложных отрицательных срабатываний — с 30 до 28. На первый взгляд это незначительное улучшение, однако важное, учитывая, что точность модели и так уже была довольно высокой. Необходимо быть осторожными при настройке модели, чтобы не настроить ее слишком сильно, так как важно сохранять баланс между чрезмерной агрессивностью и излишней

пассивностью при фильтрации спама. Пользователи предпочли бы, чтобы фильтр скорее пропускал небольшое количество спамовых сообщений, а не чтобы обычные сообщения фильтровались слишком агрессивно и расценивались как спам.

Резюме

В этой главе мы изучили классификацию с использованием наивного байесовского алгоритма. Этот алгоритм создает таблицы правдоподобия, которые затем используются для оценки вероятности того, что новые примеры относятся к тем или иным классам. Вероятности вычисляются по формуле, используемой в теореме Байеса, которая указывает на то, как связаны зависимые события. Теорема Байеса может быть слишком затратной в вычислительном отношении, однако ее упрощенная версия, которая делает так называемые наивные предположения о независимости признаков, способна обрабатывать большие наборы данных.

Наивный байесовский классификатор часто используется для классификации текста. Чтобы убедиться в его эффективности, мы использовали наивный байесовский алгоритм для задачи классификации спама в СМС. Подготовка текстовых данных для анализа потребовала применения специализированных R-пакетов для обработки и визуализации текста. В итоге модель смогла правильно классифицировать более 97 % всех СМС как спам или не спам.

Следующая глава посвящена еще двум методам машинного обучения. Каждый выполняет классификацию, разбивая данные на группы схожих значений.

5. Разделяй и властвуй: классификация с использованием деревьев решений и правил

Выбирая между предложениями о работе с разными уровнями оплаты и бонусами, многие взвешивают все за и против, руководствуясь простыми правилами, например: «Если придется добираться на работу больше часа, мне это не понравится» или «Если я заработаю менее 50 тысяч долларов, то не смогу содержать семью». Как видите, сложно сделать правильный выбор, но в этом может помочь несколько простых решений.

В этой главе будут рассмотрены деревья решений и алгоритмы обучения на основе правил — два метода машинного обучения, которые позволяют принимать сложные решения, выбирая из множества простых вариантов. Эти методы представляют решение в виде логических схем, которые можно понять без специальных знаний. Благодаря этому аспекту такие модели особенно полезны для построения бизнес-стратегий и улучшения процессов.

В этой главе вы:

- узнаете, как деревья и правила «жадно» делят данные на интересные сегменты;
- познакомитесь с самыми распространенными обучающими алгоритмами на основе дерева решений и правил классификации, включая алгоритмы C5.0, 1R и RIPPER;
- научитесь использовать эти алгоритмы для выполнения реальных задач классификации, таких как выявление рискованных банковских кредитов и распознавание ядовитых грибов.

Начнем с изучения деревьев решений, а затем перейдем к правилам классификации. После этого подведем итоги и рассмотрим краткий обзор последующих глав, в которых приведены технологии, использующие деревья и правила в качестве основы для более расширенных методов машинного обучения.

Деревья решений

Обучающие алгоритмы на базе дерева решений — это мощные классификаторы, в которых используется *древовидная структура* для моделирования отношений между признаками и возможными результатами. Как показано на рис. 5.1, эта структура получила такое название благодаря своему сходству с реальным деревом: начинается с широкого ствола, который разделяется на ветви — чем выше, тем тоньше. Примерно таким же образом в классификаторе, построенном на основе дерева решений, используется структура ветвящихся решений, по которой примеры перенаправляются к окончательному спрогнозированному значению класса.

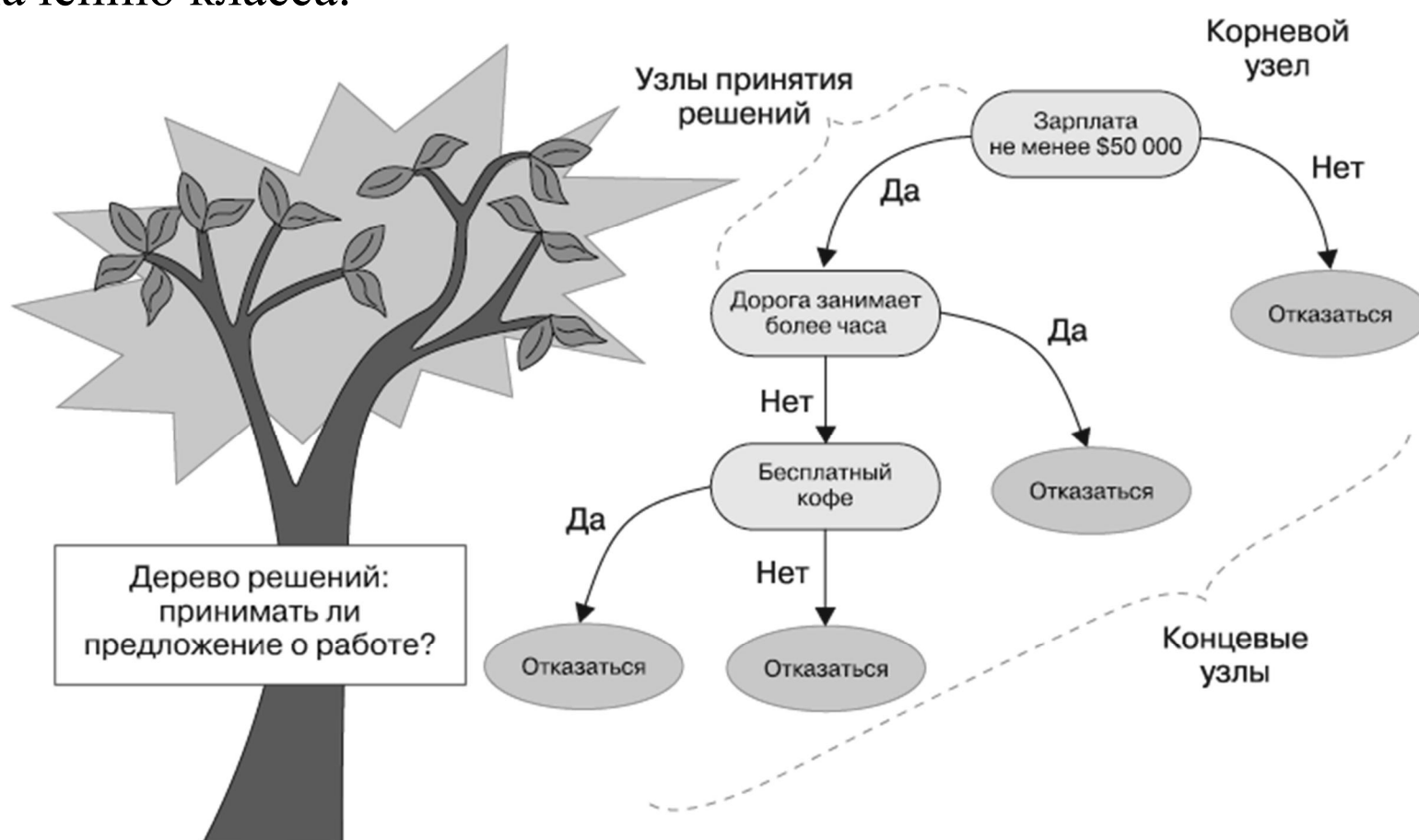


Рис. 5.1. Дерево решений, изображающее процесс принятия предложения о работе

Чтобы лучше понять, как это работает, рассмотрим дерево, которое позволяет решить, следует ли принимать предложение о работе. Рассмотрение предложения начинается с *корневого узла*; затем оно проходит через *узлы принятия решений*, в каждом из которых делается выбор на основе

условий задачи. Эти точки выбора делят данные на *ветви*, которые указывают на возможные окончательные решения. В данном случае эти решения изображены как ответы «да» и «нет», но в других случаях возможно более двух вариантов.

В точке, где может быть принято окончательное решение, дерево заканчивается *концевыми узлами* (также известны как *терминальные узлы*). Эти узлы обозначают действие, которое должно быть выполнено в результате последовательности решений. В случае модели прогнозирования концевые узлы предоставляют ожидаемый результат для данной последовательности событий в дереве.

Огромное преимущество алгоритмов, основанных на дереве принятия решений, состоит в том, что древовидная структура, подобная блок-схеме, предназначена не только для внутреннего использования машиной. После создания модели многие алгоритмы, основанные на дереве принятия решений, выводят полученную структуру в формате, удобном для восприятия человеком. Это дает представление о том, как и почему работает или не работает данная модель для конкретной задачи. Благодаря этому деревья решений особенно хорошо подходят для тех случаев, когда механизм классификации должен быть прозрачным по юридическим причинам или если результаты должны быть переданы в другие инстанции для использования в будущем. С учетом этого данные алгоритмы применяются:

- для составления рейтинга кредитоспособности, в котором причины отказа заявителю в кредите должны быть четко документированы и свободны от предвзятости;
- в маркетинговых исследованиях поведения клиентов, такого как удовлетворенность или отток, которые будут переданы руководству или рекламным агентствам;
- для диагностики заболеваний на основе лабораторных исследований, симптомов или скорости прогрессирования заболевания.

Хотя указанные области применения достаточно хорошо демонстрируют ценность деревьев для процессов обоснованного принятия решений, это не означает, что их полезность этим ограничивается. На практике деревья решений являются, пожалуй, наиболее широко используемой технологией машинного обучения и могут применяться для моделирования практически любого типа данных — часто с превосходной эффективностью без дополнительной настройки.

Тем не менее, невзирая на их широкое применение, стоит отметить, что существуют случаи, когда деревья решений далеко не идеальный вариант. Сюда входят задачи, в которых данные имеют много именованных признаков с большим количеством уровней или числовых признаков. Эти случаи могут приводить к огромному количеству решений и чрезмерно

сложным деревьям. Они также иногда способствуют склонности деревьев решений к сверхчувствительности данных, хотя, как мы скоро увидим, даже этот недостаток можно преодолеть, введя некоторые простые параметры.

Разделяй и властвуй

Деревья решений строятся с использованием эвристики, называемой *рекурсивным сегментированием*. Этот подход также широко известен как метод «разделяй и властвуй» (divide and conquer), так как он разбивает данные на подмножества, которые затем снова разделяются на еще меньшие подмножества и так далее, до тех пор, пока процесс не остановится. Это произойдет, когда алгоритм определит, что данные в подмножествах являются достаточно однородными, или выполнится другое условие остановки.

Чтобы проследить, как при разделении набора данных строится дерево решений, представьте один только корневой узел, который затем вырастает в целое дерево. Сначала корневой узел представляет собой весь набор данных, поскольку расщепления не произошло. В этот момент алгоритм построения дерева должен выбрать признак разделения; в идеале он выбирает наиболее легко предсказуемый признак целевого класса. Затем примеры (объекты) разбиваются на группы в соответствии с различными значениями этого признака и формируется первое множество ветвей дерева.

Продвигаясь по каждой ветви, алгоритм продолжает разделять и властвовать над данными, всякий раз выбирая наилучший признак-кандидат для создания очередного узла принятия решения, пока не будет достигнут критерий остановки. Метод «разделяй и властвуй» может остановиться на узле, если:

- все (или почти все) примеры в данном узле относятся к одному классу;
- не осталось признаков, по которым можно было бы разделить оставшиеся примеры;
- размер дерева достиг заданной предельной величины.

Для того чтобы лучше разобраться в процессе построения дерева, рассмотрим простой пример. Представьте, что вы работаете на голливудской киностудии. Ваша задача — решить, должна ли студия продолжать снимать фильмы по сценариям, предлагаемым молодыми многообещающими авторами. Вы вернулись из отпуска, и ваш стол завален предложениями. Чтобы не тратить время на чтение каждого предложения от начала до конца, вы решаете создать алгоритм построения дерева решений, который бы позволил спрогнозировать, попадет ли

потенциальный фильм в одну из следующих трех категорий: «Успех у критиков», «Успех в прокате» или «Полный провал».

Чтобы построить дерево решений, вы обращаетесь к архивам студии с целью изучить причины успеха и провала 30 последних фильмов компании. Вы быстро замечаете связь между бюджетом фильма, количеством знаменитостей, снявшихся в главных ролях, и успешностью фильма. Придя в восторг от этого открытия, вы строите диаграмму рассеяния (рис. 5.2).

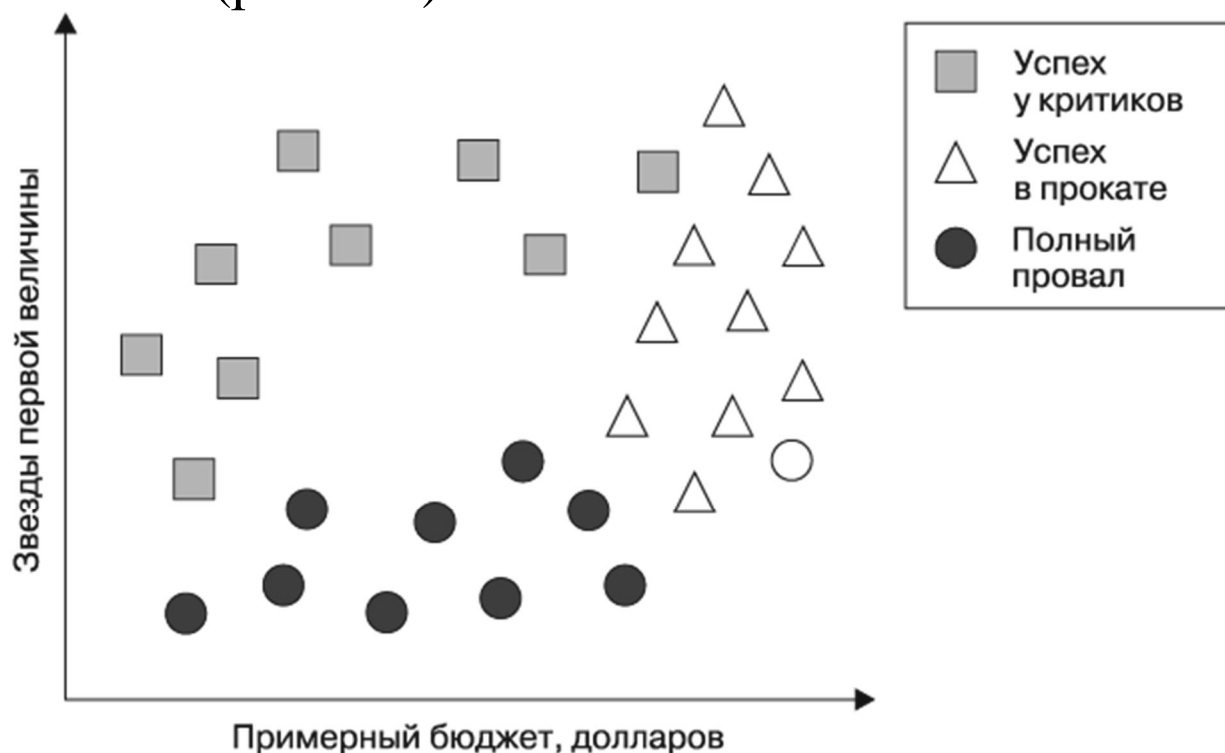


Рис. 5.2. Диаграмма рассеяния, отражающая взаимосвязь между бюджетом фильма и количеством снявшихся в нем знаменитостей

Используя стратегию «разделяй и властвуй», на основе данных можно построить простое дерево решений. Прежде всего, чтобы создать корневой узел дерева, разделим признак, указывающий на количество знаменитостей, разбивая фильмы на группы: те, где снималось много звезд первой величины, и те, где знаменитостей было мало (рис. 5.3).

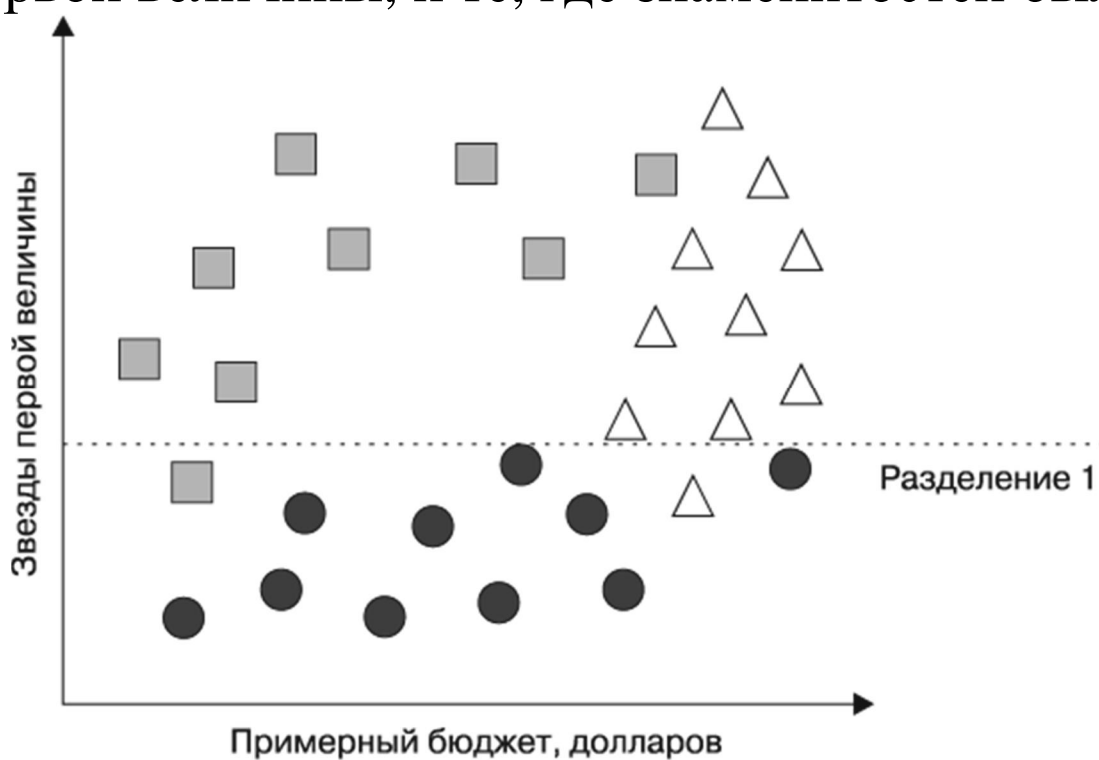


Рис. 5.3. Первое разделение дерева решений: фильмы разделены на те, в которых снималось много и мало знаменитостей

Затем группу фильмов с большим количеством знаменитостей можно разделить на еще одну: фильмы с большим и маленьким бюджетом (рис. 5.4).

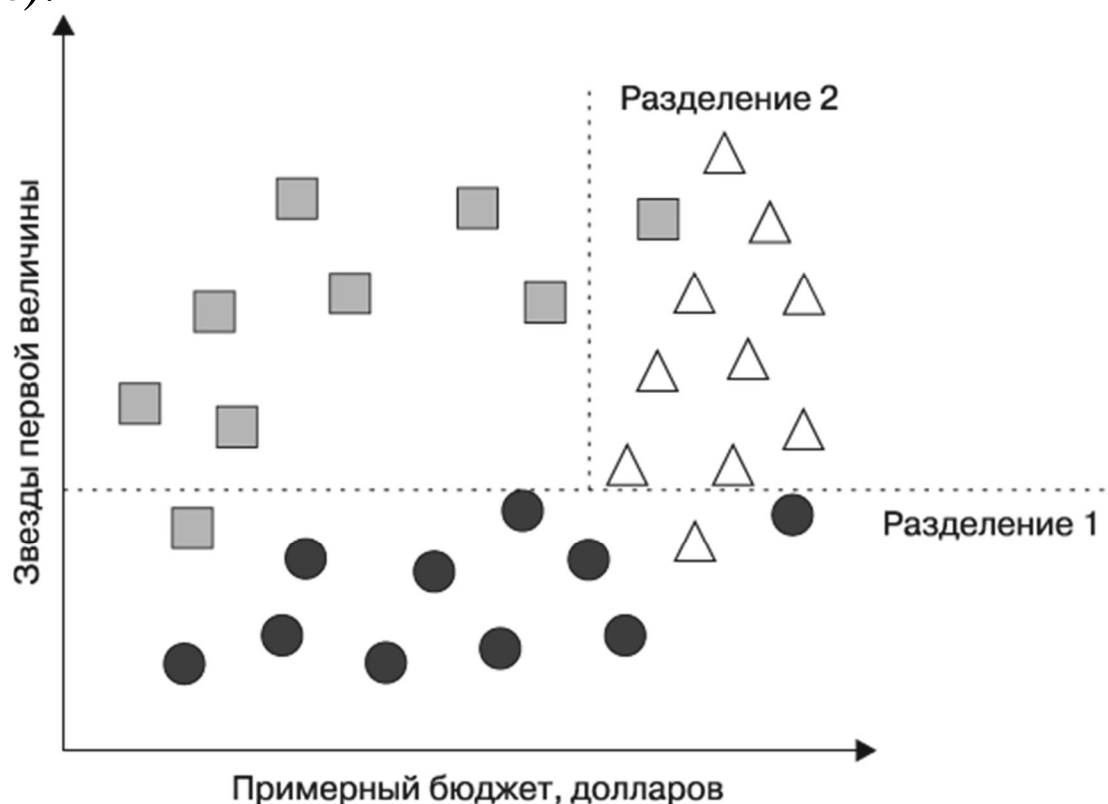


Рис. 5.4. Второе разделение дерева решений: фильмы с большим количеством знаменитостей делятся на фильмы с большим и маленьким бюджетом

В этой точке мы разделили данные на три группы. Группа, попадающая в верхний левый угол диаграммы, полностью состоит из фильмов, получивших признание у критиков. Эта группа отличается большим количеством знаменитостей и относительно низким бюджетом. В верхнем правом углу большинство фильмов — кассовые хиты с большим бюджетом и большим количеством знаменитостей. К последней группе относятся провальные фильмы, в которых снималось немного звезд, но бюджеты которых варьируются от маленьких до больших.

Если бы мы захотели, то могли бы продолжить разделять данные и властвовать над ними, разбивая их на все более специфичные диапазоны по бюджету и количеству знаменитостей, пока все неправильно классифицированные значения не будут классифицированы правильно в своем маленьком разделе. Однако не рекомендуется делать дерево решений таким сверхчувствительным. Конечно, ничто не мешает алгоритму разбивать данные до бесконечности, однако слишком специфичные решения не всегда позволяют делать более широкие обобщения. Чтобы избежать переобучения, мы остановим алгоритм на этом уровне, так как более 80 % примеров в каждой группе принадлежат к одному классу. Это будет основой для критерия остановки алгоритма.



Возможно, вы заметили, что диагональные линии могли бы разделить данные еще более четко. Но это одно из ограничений представления знаний в виде дерева решений: здесь используются только *разбиения*, *параллельные осям координат*. Тот факт, что при каждом разбиении учитывается

только один признак, не позволяет дереву формировать более сложные границы решений. Например, диагональная линия могла бы соответствовать следующему решению: «Превышает ли количество знаменитостей предполагаемый бюджет?» И если да, то результатом будет «успех у критиков».

Модель для прогнозирования успешности фильмов может быть представлена в виде простого дерева, как показано на рис. 5.5.

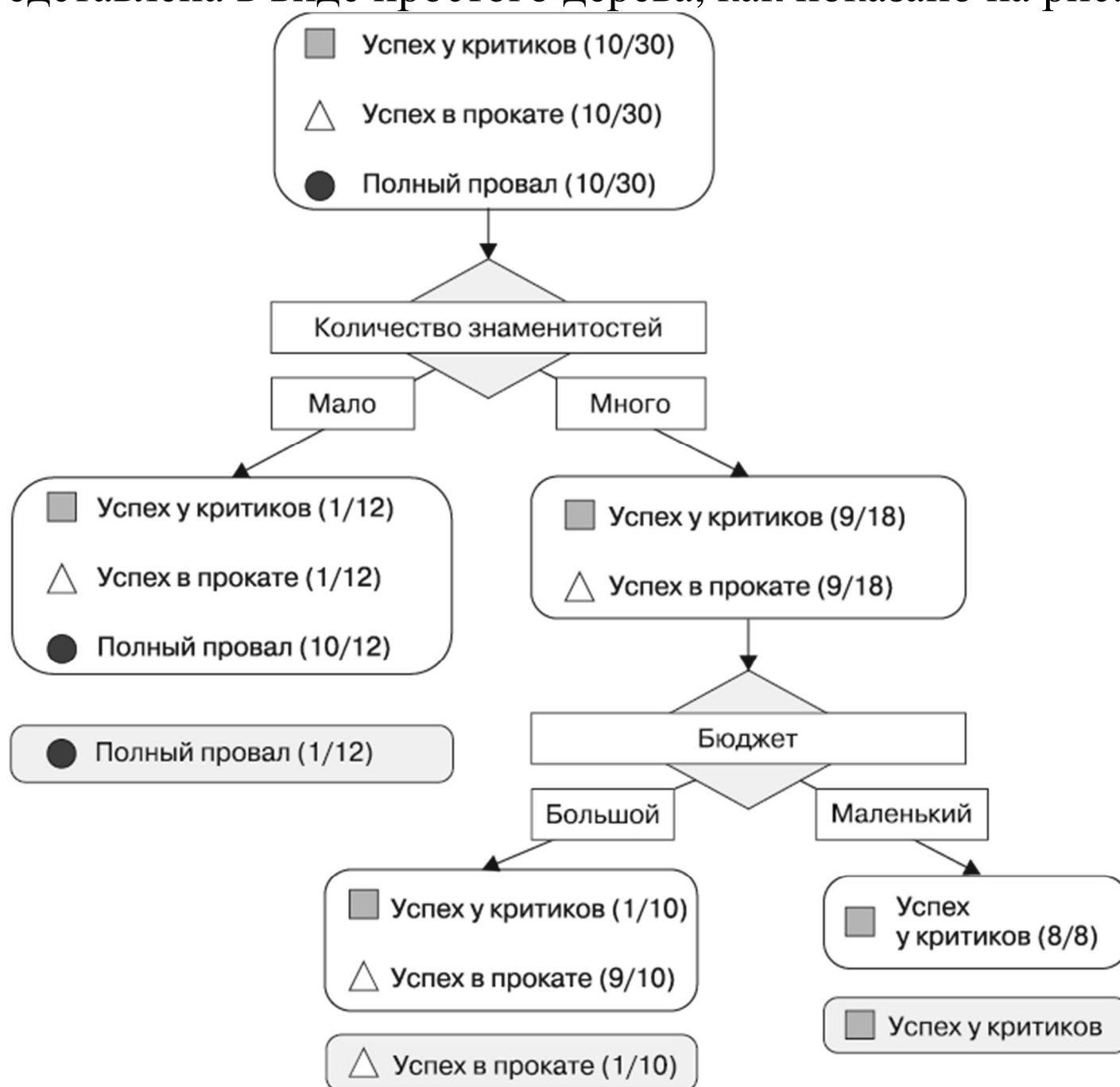


Рис. 5.5. Дерево решений, построенное на данных об уже снятых фильмах, позволяет спрогнозировать успешность будущих фильмов

Каждый шаг в дереве показывает, какая доля примеров попадает в соответствующий класс. На каждом этапе видно, как данные становятся все более однородными по мере приближения ветвей к конечным узлам. Для того чтобы оценить новый сценарий фильма, достаточно следовать указаниям по каждому решению, пока не будет спрогнозирована успешность или неудача сценария. Используя этот подход, можно быстро определить наиболее перспективные варианты среди всех сценариев и вернуться к более важной работе, такой как написание речи для выступления на церемонии вручения «Оскара».

Поскольку реальные данные содержат более двух признаков, деревья решений быстро становятся намного сложнее, чем с большим количеством ветвей, промежуточных и конечных узлов. В следующем разделе мы

познакомимся с популярным алгоритмом автоматического построения моделей дерева решений.

Алгоритм дерева решений C5.0

Существует множество реализаций деревьев решений. Одним из самых известных является *алгоритм C5.0*. Этот алгоритм разработан специалистом по компьютерным системам Дж. Россом Квинланом (J. Ross Quinlan). Это усовершенствованная версия его предыдущего алгоритма C4.5, который, в свою очередь, является улучшенным вариантом его же алгоритма *Iterative Dichotomiser 3 (ID3)*. Квинлан продает C5.0 коммерческим клиентам (подробности вы найдете на сайте <http://www.rulequest.com/>), однако исходный код однопоточной версии этого алгоритма опубликован и поэтому включен в такой программный продукт, как R.



Чтобы еще больше сбить вас с толку, в R-пакет RWeka входит J48 — популярная альтернатива C4.5 с открытым исходным кодом на основе Java. Поскольку различия между C5.0, C4.5 и J48 незначительны, принципы, изложенные в этой главе, применимы к любому из этих трех методов и алгоритмы следует считать синонимичными.

Алгоритм C5.0 стал отраслевым стандартом для построения деревьев решений, поскольку он без дополнительной настройки хорошо подходит для решения большинства типов задач. По сравнению с другими расширенными моделями машинного обучения, такими как описанные в главе 7, деревья решений, построенные на базе C5.0, обычно работают почти так же хорошо, но их гораздо легче понять и реализовать. Кроме того, как показано в табл. 5.1, недостатки этого алгоритма сравнительно невелики и их по большей части можно избежать.

Таблица 5.1

Преимущества	Недостатки
Универсальный классификатор, который хорошо справляется со многими типами задач. Высокоавтоматизированный процесс обучения, позволяющий обрабатывать числовые и номинальные характеристики, а также отсутствующие данные. Исключает неважные признаки. Может использоваться и для маленьких, и для больших наборов данных. Приводит к построению модели, которая может быть интерпретирована без математического образования (для относительно небольших деревьев). Эта модель эффективнее, чем другие, более сложные модели	Модели дерева решений часто смещены в сторону разделения по признакам, имеющим большое количество уровней. Легко получить как переобученную, так и недостаточно обученную модель. При моделировании некоторых отношений возможны проблемы из-за разделений только вдоль осей координат. Незначительные изменения тренировочных данных могут привести к значительным изменениям в логике принятия решений. Большие деревья бывает трудно интерпретировать, а решения, которые они рекомендуют, могут показаться нелогичными

Для простоты в предыдущем примере дерева решений игнорировалась математика, связанная с тем, как машина использует стратегию «разделяй

и властью». Рассмотрим это более подробно, чтобы изучить, как эвристика работает на практике.

Выбор лучшего разделения

Первая проблема, с которой сталкивается дерево решений, — определение того, по какому признаку осуществлять разделение. В предыдущем примере мы искали способ разделения данных таким образом, чтобы в полученных разделах содержались примеры, главным образом относящиеся к одному классу. Степень, в которой подмножество примеров содержит только один класс, называется *чистотой*, а подмножество, состоящее только из одного класса, называется *чистым*.

Существуют различные показатели чистоты, которые можно использовать для определения наилучшего кандидата для разделения дерева решений. В C5.0 используется *энтропия* — концепция, заимствованная из теории информации, которая количественно определяет случайность, или беспорядочность, множества значений класса. Множества с высокой энтропией очень разнообразны и предоставляют мало информации о других элементах, которые также могут принадлежать к этому множеству, поскольку между ними нет очевидной общности. Дерево решений старается найти такие разделения, которые уменьшают энтропию, в итоге повышая однородность внутри групп.

Обычно энтропия измеряется в *битах*. Если существует только два возможных класса, то энтропия принимает значения 0 и 1. Для n классов энтропия принимает значения от 0 до $\log_2(n)$. В каждом случае минимальное значение указывает на то, что выборка является полностью однородной, а максимальное — на то, что данные настолько разнообразны, насколько это возможно, и ни у одной группы нет ничего общего.

В математическом представлении энтропия (Entropy) определяется следующим образом:

$$\text{Entropy}(s) = \sum_{i=1}^c -p_i \log_2(p_i)$$

В этой формуле для выбранного сегмента данных $(s)_c$ означает число уровней классов, а p_i — долю значений, попадающих в класс уровня i . Например, предположим, что у нас есть раздел данных с двумя классами: красный (60 %) и белый (40 %). Тогда энтропию можно рассчитать следующим образом:

```
> -0.60 * log2(0.60) - 0.40 * log2(0.40) [1]  
0.9709506
```

Мы можем визуализировать энтропию для всех возможных вариантов, состоящих из двух классов. Если известно, что доля примеров, относящихся к одному классу, равна x , то доля примеров, принадлежащих

к другому классу, равна $(1 - x)$. Используя функцию `curve()`, можно построить график энтропии для всех возможных значений x :

```
> curve(-x * log2(x) - (1 - x) * log2(1 - x), col = "red", xlab = "x", ylab = "Entropy", lwd = 4)
```

В результате получим следующий график (рис. 5.6).

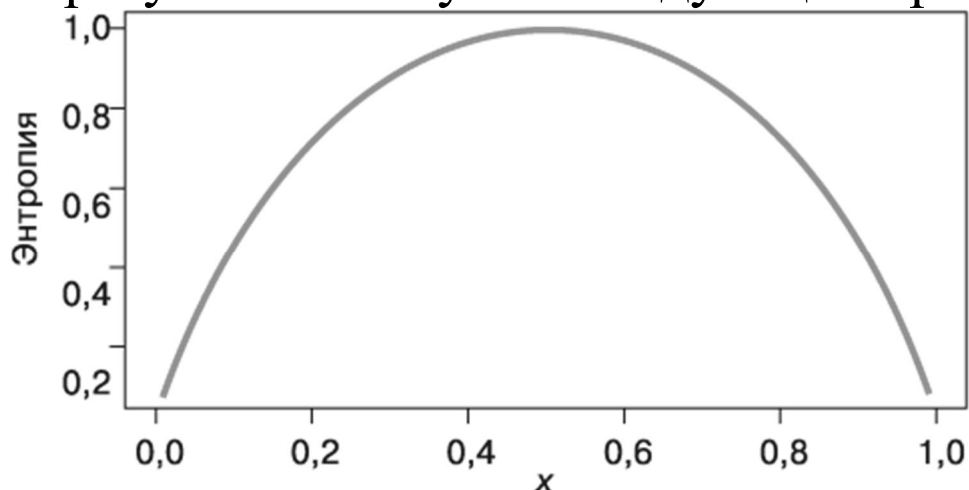


Рис. 5.6. Общая энтропия как пропорция изменений одного класса к результату двух классов

Как видно на графике, энтропия достигает пикового значения при $x = 0,50$: разделение на уровне 50–50 приводит к максимальной энтропии. Поскольку один класс все больше доминирует в другом, энтропия снижается до нуля.

Чтобы использовать энтропию для определения оптимального признака разделения, алгоритм рассчитывает изменение однородности, которое может возникнуть при разделении по всем возможным признакам. Эта мера называется *приростом информации*. Прирост информации (infogain) для признака F вычисляется как разность между энтропией в данном сегменте до разделения (s_1) и в сегментах, полученных в результате разделения (s_2):

$$\text{InfoGain}(F) = \text{Entropy}(s_1) - \text{Entropy}(s_2).$$

Сложность этого подхода заключается в том, что в результате разделения создается несколько сегментов данных. Таким образом, функция вычисления энтропии $\text{Entropy}(s_2)$ должна учитывать общую энтропию для всех сегментов. Ее можно получить путем взвешивания энтропии каждого сегмента в соответствии с долей всех записей, попадающих в этот сегмент. Это можно выразить в формуле:

$$\text{Entropy}(s) = \sum_{i=1}^n w_i \text{Entropy}(P_i).$$

Проще говоря, общая энтропия, полученная в результате разделения, представляет собой сумму энтропии каждого из n сегментов, взвешенную в соответствии с долей примеров, попадающих в этот сегмент (w_i).

Чем больше прирост информации, тем лучше данный признак подходит для создания однородных групп после разделения по этому признаку. Если прирост информации равен нулю, то при разделении по этому признаку энтропия не уменьшается. Однако максимальный прирост информации

равен энтропии до разделения. Таким образом, если энтропия после разделения равна нулю, это означает, что разделение привело к созданию полностью однородных групп.

Приведенные выше формулы рассчитаны на номинальные признаки, однако деревья решений также позволяют использовать прирост информации для разделения по числовым признакам. Для этого обычно проверяют различные варианты разделений, при которых значения делятся на группы, где значения данного признака больше или меньше некоего порогового значения. Это позволяет свести числовой признак к двухуровневому категориальному признаку, после чего вычислять прирост информации как обычно. Для разделения выбирается числовая точка среза, обеспечивающая максимальный прирост информации.



Прирост информации, используемый в C5.0, не единственный критерий разделения, который можно использовать для построения деревьев решений. Другими часто используемыми критериями являются индекс Джини, критерий хи-квадрата и коэффициент усиления. Обзор этих (и многих других) критериев вы найдете в статье: Mingers J. An Empirical Comparison of Selection Measures for Decision-Tree Induction. Machine Learning, 1989. Vol. 3.P. 319–342.

Сокращение дерева решений

Как уже отмечалось, дерево решений может расти бесконечно, выбирая признаки разделения и ветвясь на все меньшие сегменты, пока не будут полностью классифицированы все примеры или пока в алгоритме не исчерпаются возможности для разделения. Но, когда дерево становится слишком большим, многие принимаемые им решения оказываются чрезмерно конкретными, и модель становится переобученной.

Процесс *сокращения* дерева решений означает уменьшение его размера таким образом, чтобы дерево лучше обобщало новые данные.

Одним из решений этой проблемы является остановка роста дерева при достижении определенного числа точек принятия решений или когда узлы принятия решений будут содержать лишь небольшое количество примеров. Это называется *ранней остановкой* ИЛИ *ранним сокращением* дерева решений. Поскольку при построении дерева алгоритм стремится избегать ненужной работы, это привлекательная стратегия. Одним из недостатков такого подхода является то, что не существует способа узнать, не будут ли при этом пропущены неявные, но важные паттерны, которые были бы изучены, если бы дерево выросло до большего размера.

Альтернатива, называемая *поздним сокращением*, подразумевает построение намеренно слишком большого дерева и сокращение концевых узлов, чтобы

уменьшить размер дерева до приемлемого. Такой подход часто оказывается более эффективным, чем ранняя остановка, поскольку бывает довольно сложно определить оптимальную глубину дерева решений заранее, без предварительного увеличения. Последующее сокращение дерева позволяет алгоритму гарантировать, что были обнаружены все важные структуры данных.



Подробности реализации операций сокращения дерева очень специфичны и выходят за рамки этой книги. Сравнение некоторых существующих методов вы найдете в статье: Esposito F., Malerba D., Semeraro G. A Comparative Analysis of Methods for Pruning Decision Trees // IEEE Transactions on Pattern Analysis and Machine Intelligence, 1997. Vol. 19. P. 476–491.

Одним из преимуществ алгоритма C5.0 является то, что сокращение в нем сильно автоматизировано: алгоритм принимает многие решения самостоятельно, используя весьма разумные значения, предлагаемые по умолчанию. Общая стратегия этого алгоритма заключается в позднем сокращении дерева. Сначала строится большое дерево, соответствующее тренировочным данным. Затем узлы и ветви, которые мало влияют на ошибки классификации, удаляются. В некоторых случаях целые ветви перемещаются вверх по дереву или заменяются более простыми решениями. Эти процессы пересадки ветвей называются *поднятием поддерева* и *заменой поддерева* соответственно.

Достижение правильного баланса между переобучением и недообучением дерева — это в некотором роде искусство. Однако, если точность модели жизненно важна, возможно, стоит потратить некоторое время на исследование различных вариантов сокращения, чтобы посмотреть, не повысится ли от этого эффективность распознавания тестового набора данных. Как вы вскоре увидите, одним из преимуществ алгоритма C5.0 является то, что он очень легко позволяет настроить параметры обучения.

Пример: распознавание рискованных банковских кредитов с помощью деревьев решений C5.0

Мировой финансовый кризис 2007–2008 годов показал, как важна прозрачность и строгость в принятии банковских решений. Когда кредиты стали менее доступными, банки ужесточили систему кредитования и обратились к машинному обучению для более точного определения рискованных кредитов.

Благодаря высокой точности и возможности формулировать статистическую модель на понятном человеку языке дерева решений

широко применяются в банковской сфере. Поскольку правительства многих стран тщательно следят за справедливостью кредитования, руководители банков должны быть в состоянии объяснить, почему одному заявителю было отказано в получении займа, в то время как другому одобрили выдачу кредита. Эта информация полезна и для клиентов, желающих узнать, почему их кредитный рейтинг оказался неудовлетворительным.

Похоже, что автоматические модели оценки кредитоспособности используются для рассылок по кредитным картам и мгновенных онлайн-процессов одобрения кредитов. В этом разделе мы разработаем простую модель принятия решения о предоставлении кредита с использованием алгоритма построения деревьев решений C5.0. Вы также узнаете, как настроить параметры модели, чтобы свести к минимуму ошибки, которые могут привести к финансовым потерям.

Шаг 1. Сбор данных

Цель кредитной модели — выявить факторы, связанные с более высоким риском невозвращения кредита. Для этого необходимо получить данные о большом количестве предыдущих банковских кредитов и информацию о получателях этих кредитов.

Данные с такими характеристиками доступны в наборе данных, переданном Хансом Хофманном (Hans Hofmann) из Гамбургского университета в репозиторий машинного обучения UCI (<http://archive.ics.uci.edu/ml>). Этот набор данных содержит информацию о кредитах, полученных от кредитного агентства в Германии.



Набор данных, представленный в этой главе, был немного изменен по сравнению с оригиналом, чтобы исключить некоторые этапы предварительной обработки. Чтобы выполнить эти примеры, загрузите файл `credit.csv` и сохраните его в рабочем каталоге среды R. Набор кредитных данных включает в себя 1000 примеров предоставления кредитов, а также набор числовых и номинальных признаков, определяющих характеристики кредита и заявителя. Переменная `class` указывает на то, был ли возвращен кредит. Посмотрим, можно ли выделить какие-либо модели, которые бы прогнозировали этот результат.

Шаг 2. Исследование и подготовка данных

Как и ранее, импортируем данные, используя функцию `read.csv()`. Поскольку символьные данные являются полностью категоризованными, можно пропустить параметр `stringsAsFactors`, который по

умолчанию принимает значение `TRUE`. В результате будет создан фрейм данных о кредитах с несколькими факторными переменными:

```
> credit <- read.csv("credit.csv")
```

Проверить полученный объект можно, просмотрев первые несколько строк вывода функции `str()`:

```
> str(credit)'data.frame':1000 obs. of 17
variables: $ checking_balance : Factor w/ 4
levels "< 0 DM", "> 200 DM", .. $
months_loan_duration: int 6 48 12 ... $
credit_history : Factor w/ 5 levels
"critical", "good", .. $ purpose
Factor w/ 6 levels "business", "car", .. $
amount : int 1169 5951 2096 ...
```

Как и ожидалось, здесь видим 1000 наблюдений и 17 признаков, часть из которых относится к факторному, а часть — к целочисленному типу данных.

Посмотрим на вывод функции `table()` для нескольких признаков кредита, которые, очевидно, прогнозируют невозвращение кредита. Баланс текущего и накопительного счетов заявителя представлен в виде категориальных переменных:

```
> table(credit$checking_balance) < 0 DM >
200 DM 1 - 200
DM unknown 274 63
269 394>
table(credit$savings_balance) < 100 DM >
1000 DM 100 - 500 DM 500 - 1000
DM unknown 603 48
103 63 183
```

Баланс на текущем и накопительном счетах может быть важным фактором, определяющим вероятность невозвращения кредита. Обратите внимание, что, поскольку данные о кредите были получены из Германии, то значения измеряются в немецких марках (DM) — валюте, которая использовалась в Германии до введения евро.

Некоторые из признаков кредита являются числовыми — например, продолжительность и запрашиваемая сумма:

```
>
summary(credit$months_loan_duration) Min.
1st Qu. Median Mean 3rd
Qu. Max. 4.0 12.0 18.0 20.9
24.0 72.0>
```

summary(credit\$amount)		Min.	1st	2nd	3rd	Max.
Qu.	Median	Mean	3rd			
Qu.	Max.	250	1366	2320	3271	
		3972	18420			

Суммы кредита варьировались от 250 до 18 420 DM и предоставлялись на срок от 4 до 72 месяцев. Средняя сумма составила 2320 DM, а средняя продолжительность — 18 месяцев.

Вектор `default` указывает на то, смог ли получатель кредита выполнить установленные условия платежа, или же он обанкротился. В общей сложности 30 % кредитов из этого набора данных закончились банкротством:

```
> table(credit$default)no    yes700    300
```

Высокий уровень банкротств нежелателен для банка, поскольку это означает, что банк едва ли сможет окупить инвестиции. Если мы решим задачу, то модель сможет определять соискателей, имеющих высокий риск банкротства, что позволит банку отклонять запросы на кредит, вместо того чтобы выдавать деньги.

Подготовка данных: создание случайных обучающего и тестового наборов данных

Как и в предыдущих главах, разделим данные на две части: тренировочный набор для построения дерева решений и тестовый набор для оценки эффективности дерева решения на новых данных. Используем 90 % данных для обучения и 10 % — для тестирования, что даст нам 100 записей для моделирования новых соискателей.

В предыдущих главах использовались данные, которые были заранее представлены в случайном порядке, и поэтому мы просто разделяли набор данных на две части, взяв первые 90 % записей для обучения и оставшиеся 10 % — для тестирования. Однако набор кредитных данных не упорядочен случайным образом, поэтому такой подход был бы неразумным. Предположим, что банк отсортировал данные по сумме кредита, причем кредиты на самую большую сумму оказались в конце файла. Если бы мы использовали первые 90 % этих данных для обучения и оставшиеся 10 % для тестирования, то мы бы обучали модель только на небольших кредитах, а тестировали ее — на кредитах на большие суммы. Очевидно, это может быть проблематично.

Решим эту проблему, обучив модель на *случайной выборке* кредитных данных. Случайная выборка — это процесс, который случайным образом выбирает подмножество записей. В R для выполнения случайной выборки используется функция `sample()`. Однако, прежде чем приступить к работе, обычной практикой является выбор *начального значения*, которое заставляет процесс рандомизации выполнять последовательность, которая

позже может быть воспроизведена. Может показаться, что это делает бессмысленным генерирование случайных чисел, но для этого есть веская причина. Предоставление начального значения с помощью функции `set.seed()` гарантирует, что если потом анализ будет повторяться, то будет получен такой же результат.



Возможно, вы спросите: как задать начальное значение для так называемого случайного процесса, чтобы каждый раз получать одинаковые результаты? Дело в том, что для создания последовательностей случайных чисел в компьютерах используется математическая функция, называемая генератором псевдослучайных чисел. Эти числа очень похожи на случайные, однако на самом деле вполне предсказуемы, если знать предыдущие значения последовательности. На практике современные последовательности псевдослучайных чисел практически неотличимы от истинных случайных последовательностей, но имеют преимущество — компьютеры быстро и легко их генерируют.

В следующих командах используется функция `sample()` с начальным значением. Обратите внимание, что функция `set.seed()` использует произвольное значение `123`. Если не указать начальное значение, то разделение данных на тренировочный и тестовый наборы будет отличаться от того, что показано дальше в этой главе. Следующие команды случайным образом выбирают 900 значений из последовательности целых чисел от 1 до 1000:

```
> set.seed(123) > train_sample <- sample(1000, 900)
```

Как и ожидалось, полученный в результате объект `train_sample` представляет собой вектор из 900 случайных целых чисел:

```
> str(train_sample) int [1:900] 288 788 409 881 937 46 525 887 548 453 ...
```

Используя этот вектор для выбора строк из кредитных данных, можно разделить эти данные на 90 % тренировочных и 10 % тестовых данных, как мы и хотели. Напомню, что оператор отрицания (символ `-`), используемый при выборе тестовых записей, дает среде R команду выбирать те записи, которые не содержатся в указанных строках; другими словами, тестовые данные включают в себя только те строки, которые отсутствуют в обучающей выборке:

```
> credit_train <- credit[train_sample, ] > credit_test <- credit[-train_sample, ]
```

Если рандомизация выполнена правильно, то в каждом наборе данных должно быть примерно 30 % невозвращенных кредитов:

```
>
prop.table(table(credit_train$default))      no
      yes 0.7033333      0.2966667>
prop.table(table(credit_test$default))      no      ye
s0.67      0.33
```

Поскольку в обучающем и в тестовом наборах данных количество невозвращенных кредитов оказалось примерно одинаковым, теперь мы можем построить дерево решений. Если бы пропорции сильно различались, можно было бы принять решение о повторной выборке набора данных или попытаться применить более сложный алгоритм выборки, например такой, как описан в главе 10.



Если ваши результаты не совпадают с моими, проверьте, выполнили ли вы команду `set.seed(123)` непосредственно перед созданием вектора `train_sample`.

Шаг 3. Обучение модели на данных

Для обучения данной модели дерева решений воспользуемся алгоритмом C5.0 из пакета `C50`. Если вы этого еще не сделали, установите этот пакет с помощью команды `install.packages("C50")` и загрузите его в R-сессию, выполнив команду `library(C50)`.

В следующей таблице перечислены некоторые параметры, наиболее часто используемые при построении деревьев решений. По сравнению с описанными ранее подходами машинного обучения, алгоритм C5.0 предлагает гораздо больше способов адаптировать модель к конкретной задаче обучения.

После загрузки пакета `C50` команда `?C5.0Control` позволяет увидеть страницу справки с дополнительной информацией о том, как точнее настроить алгоритм, используя следующие параметры.

Синтаксис алгоритма дерева решений C5.0

Использование функции `c5.0()` из пакета `C50`

Построение классификатора:

```
m <- c5.0(train, class, trials = 1, costs = NULL)
```

`train` – фрейм данных, содержащий данные обучения;

`class` – факторный вектор, содержащий классы для всех строк тренировочных данных;

`trials` – дополнительный параметр: число, определяющее количество усиливающих деревьев (по умолчанию равен 1);

`costs` – дополнительный параметр: матрица, определяющая цену для различных типов ошибок.

Функция возвращает модель C5.0, которую можно использовать для прогнозирования.

Прогнозирование:

```
p <- predict(m, test, type = "class")
```

`m` – модель, обученная с помощью функции `c5.0()`;

`test` – фрейм данных, содержащий тестовые данные с теми же признаками, что и тренировочные данные, использованные для построения классификатора;

`type` – переменная, принимающая одно из двух значений, "class" или "prob", и определяющая, следует ли спрогнозировать наиболее вероятное значение класса или же необработанную вероятность.

Функция возвращает вектор спрогнозированных значений классов или необработанных вероятностей, в зависимости от значения параметра `type`.

Пример:

```
credit_model <- c5.0(credit_train, loan_default)
credit_prediction <- predict(credit_model, credit_test)
```

Для первой итерации модели предоставления кредита используем настройки C5.0, предлагаемые по умолчанию, как показано в следующем коде. В столбце 17 фрейма данных `credit_train` содержится переменная класса `default`, поэтому исключим этот столбец из фрейма данных обучения и представим его в качестве целевого фактора для классификации:

```
> credit_model <- C5.0(credit_train[-17],
credit_train$default)
```

Теперь объект `credit_model` содержит дерево решений C5.0. Чтобы увидеть некоторые основные сведения об этом дереве, достаточно ввести его имя:

```
> credit_modelCall:C5.0.default(x =
credit_train[-17], y =
credit_train$default)Classification TreeNumber of
samples: 900Number of predictors: 16Tree size:
57Non-standard options: attempt to group
attributes
```


Как видим, здесь показаны некоторые базовые факты о дереве, включая вызов функции, которая его сгенерировала, количество признаков (здесь они названы `predictors`) и примеров (называемых `samples`), которые были использованы при построении дерева. Указан и размер дерева — 57: это означает, что глубина дерева составляет 57 решений. Существенно больше, чем те примеры деревьев, которые мы рассматривали до сих пор!

Чтобы увидеть дерево решений, можно вызвать функцию `summary()` для модели:

```
> summary(credit_model)
```

Получим следующий результат:

```
C5.0 [Release 2.07 GPL Edition]
```

```
-----
```

```
Class specified by attribute `outcome`
```

```
Read 900 cases (17 attributes) from undefined.data
```

```
Decision tree:
```

```
checking_balance in {> 200 DM,unknown}: no (412/50)
```

```
checking_balance in {< 0 DM,1 - 200 DM}:
```

```
...credit_history in {perfect,very good}: yes (59/18)
```

```
  credit_history in {critical,good,poor}:
```

```
    ...months_loan_duration <= 22:
```

```
      ...credit_history = critical: no (72/14)
```

```
      : credit_history = poor:
```

```
      :   ...dependents > 1: no (5)
```

```
      :   : dependents <= 1:
```

```
      :   :   ...years_at_residence <= 3: yes (4/1)
```

```
      :   :   years_at_residence > 3: no (5/1)
```

Здесь показаны несколько первых ветвей дерева решений. Первые три строки в переводе на обычный язык означают следующее.

- Если остаток на текущем счете неизвестен или превышает 200 марок, кредит классифицируется как «маловероятно, что он не будет возвращен».
- В противном случае, если остаток на текущем счете отрицательный или меньше 200 марок...
- ...и у заявителя отличная или очень хорошая кредитная история, то кредит классифицируется как «с высокой вероятностью невозвращения».

Цифры в скобках указывают на количество примеров, удовлетворяющих критериям этого решения, и число ошибочно классифицированных решений. Например, первая строка с цифрами 412/50 означает, что из 412 примеров, удовлетворяющих этому решению, 50 были ошибочно классифицированы как «маловероятно, что не будет возвращен». Другими словами, 50 заявителей на самом деле не вернули кредит, несмотря на то что модель спрогнозировала обратное.



Иногда дерево приводит к решениям, которые являются бессмысленными с точки зрения логики. Например, почему заявитель, имеющий очень хорошую кредитную историю, может не выполнить своих обязательств, в то время как те, чей баланс на текущем счете неизвестен, классифицируются как несклонные к банкротству? Иногда алгоритм создает подобные противоречивые правила. Они могут отражать как реальную закономерность данных, так и статистическую аномалию. В любом случае важно исследовать подобные странные решения, чтобы понять, имеет ли логика дерева смысл для использования в бизнесе. После построения дерева результат (`credit_model`) отображает матрицу несоответствий — перекрестную таблицу, в которой указаны записи тренировочных данных, неправильно классифицированные моделью:

```
Evaluation on training data (900
cases):      Decision Tree      -----
      Size      Errors      56      133(14.8%)      <<
      (a)      (b)      <-classified as      ----      ----
      598      35      (a): class
no      98      169      (b): class yes
```

Заголовок `Errors` указывает на то, что модель правильно классифицировала все записи, кроме 133 из 900 тренировочных экземпляров, так что коэффициент ошибок составил 14,8 %. В итоге 35 реальных значений `no` были ошибочно классифицированы как `yes` (ложные положительные срабатывания), а 98 значений `yes` были ошибочно классифицированы как `no` (ложные отрицательные срабатывания).

Учитывая тенденцию деревьев решений к переобучению на тренировочных данных, приведенная здесь частота ошибок, полученная на основании эффективности тренировочных данных, может быть чрезмерно оптимистичной. Поэтому особенно важно продолжить оценивать алгоритм, применив дерево решений к тестовому набору данных.

Шаг 4. Оценка эффективности модели

Чтобы применить дерево решений к набору тестовых данных, используем функцию `predict()`, как показано в следующей строке кода:

```
> credit_pred <- predict(credit_model,
credit_test)
```

Эта команда создает вектор спрогнозированных значений класса, которые можно затем сравнить с реальными значениями класса, используя функцию `CrossTable()` из пакета `gmodels`. Если присвоить параметрам `prop.c` и `prop.r` значение `FALSE`, то процентное

содержание столбцов и строк будет удалено из таблицы. Оставшийся процент (`prop.t`) указывает долю записей в ячейке от общего количества записей:

```
> library(gmodels)>
CrossTable(credit_test$default,
credit_pred, prop.chisq = FALSE,
prop.c = FALSE, prop.r = FALSE, dnn =
c('actual default', 'predicted default'))
```

В результате получим следующую таблицу:

actual default	predicted default		Row Total
	no	yes	
no	59 0.590	8 0.080	67
yes	19 0.190	14 0.140	33
Column Total	78	22	100

Данная модель правильно определила, что из 100 заявлений на получение кредита, входящих в тестовый набор данных, 59 не закончились банкротством, а в 14 случаях кредит не был возвращен. Это означает, что 73 % прогнозов были сделаны точно, а количество ошибок составило 27 %. Это несколько хуже, чем эффективность алгоритма на тренировочных данных, но не является неожиданным, учитывая, что эффективность модели на незнакомых данных обычно снижается. Обратите также внимание на то, что в тестовых данных модель правильно спрогнозировала лишь 14 из 33 реальных невозвращений кредита, то есть 42 %. К сожалению, этот тип ошибок потенциально является очень дорогостоящим, поскольку при каждом невозвращении кредита банк теряет деньги. Посмотрим, можно ли улучшить результат, приложив немного больше усилий.

Шаг 5. Повышение эффективности модели

Для реализации в виде реального приложения оценки кредитоспособности вероятность появления ошибок в нашей модели слишком высока.

Фактически, если бы модель прогнозировала успешный возврат кредита для каждого тестового случая, она была бы права в 67 % случаев — результат не намного хуже, чем в нашей модели, но требует гораздо меньших усилий! Прогнозирование невозвращения кредита на 900 примеров — задача сложная.

Что еще хуже, наша модель особенно плохо работает при выявлении кандидатов, которые не выполняют своих обязательств по кредитам. К счастью, существует несколько простых способов настройки алгоритма C5.0, позволяющих повысить эффективность модели в целом.

Повышение точности алгоритма построения деревьев решений
Одно из усовершенствований алгоритма C5.0 по сравнению с C4.5 — *адаптивное усиление*. Это процесс, при котором строится несколько деревьев решений, и лучший класс для каждого примера выбирается среди этих деревьев путем голосования.



Идея усиления основана главным образом на исследованиях Роба Шапира (Rob Schapire) и Йоава Фрейнда (Yoav Freund). Для получения дополнительной информации поищите в Интернете их публикации или обратитесь к их учебнику *Boosting: Foundations and Algorithms*. Cambridge, MA, The MIT Press, 2012.

Поскольку в более широком смысле усиление может применяться к любому алгоритму машинного обучения, оно будет рассмотрено более подробно далее, в главе 11. На данный момент достаточно сказать, что усиление основано на том, что, объединяя несколько слабых моделей обучения, можно создать их группу, которая будет намного сильнее, чем каждая из моделей в отдельности. Каждая модель обладает уникальным набором сильных и слабых сторон и может работать лучше или хуже, в зависимости от конкретной задачи. Таким образом, сочетая несколько моделей обучения, у каждой из которых есть свои сильные и слабые стороны, можно значительно повысить точность классификатора.

Функция `C5.0()` позволяет легко добавить усиление в дерево решений. Достаточно указать дополнительный параметр `trials`, определяющий количество отдельных деревьев решений, которые будут использоваться в усиленной команде. Параметр `trials` устанавливает верхний предел; алгоритм прекратит добавлять деревья, если обнаружит, что дополнительные проверки не улучшают точность. Мы начнем с десяти проверок — это число стало стандартом де-факто, поскольку, как показывают исследования, такое количество проверок снижает частоту появления ошибок в тестовых данных примерно на 25 %. За исключением нового параметра, команда аналогична предыдущей:

```
> credit_boost10 <- C5.0(credit_train[-17],  
credit_train$default, trials = 10)
```

Исследуя полученную модель, мы видим, что теперь в результате указано наличие усиливающих деревьев:

```
> credit_boost10Number of boosting iterations:  
10Average tree size: 47.5
```

Новый результат показывает, что за десять итераций размер дерева сократился. Если хотите, можете просмотреть все десять деревьев, введя в

командной строке `summary(credit_boost10)`. Выходные данные также показывают эффективность дерева на тренировочных данных:

```
>
summary(credit_boost10)           (a)      (b)      <-
classified as      -----
                629      4      (a): class
no                30     237     (b): class yes
```

В 900 учебных примерах классификатор допустил 34 ошибки, так что коэффициент ошибок составил 3,8 %. Это значительное улучшение по сравнению с 13,9%-ной ошибкой обучения, которую мы получили до того, как ввели усиление! Однако еще неизвестно, получим ли мы подобное улучшение на тестовых данных. Проверим:

```
> credit_boost_pred10 <-
predict(credit_boost10, credit_test)>
CrossTable(credit_test$default,
credit_boost_pred10,
           prop.chisq =
FALSE, prop.c = FALSE, prop.r =
FALSE,
           dnn = c('actual default',
'predicted default'))
```

Полученная таблица выглядит следующим образом:

actual default	predicted default		Row Total
	no	yes	
no	62 0.620	5 0.050	67
yes	13 0.130	20 0.200	33
Column Total	75	25	100

Здесь общий коэффициент ошибок снизился с 27 % без усиления до 18 % в модели с усилением. Это может показаться не столь значительным, однако на самом деле это больше, чем ожидаемое снижение до 25 %. Однако такая модель все еще не справляется с прогнозированием невозвращения кредитов, правильно предсказывая только $20 / 33 = 61$ %. Отсутствие улучшения может быть следствием относительно небольшого тренировочного набора данных, или же просто задача может быть слишком трудной.

И все же, если усиление так легко добавить, почему бы не применять его по умолчанию к каждому дереву решений? На то есть две причины. Во-первых, если построение дерева решений занимает много времени, то построение множества деревьев может оказаться нецелесообразным в вычислительном отношении. Во-вторых, если тренировочные данные очень зашумленные, то усиление может вообще не привести к улучшению.

Тем не менее, если требуется большая точность, стоит попробовать применить усиление.

Цена ошибки бывает разной

Предоставление кредита заявителю, который, скорее всего, не в силах выполнить свои обязательства, может дорого обойтись. Одним из решений, позволяющих сократить количество ложных отрицательных результатов, может стать более жесткий отказ в пограничных случаях, исходя из предположения, что проценты, которые банк получит от рискованного кредита, и близко не окупят тех огромных потерь, которые банк понесет, если деньги не будут выплачены вообще.

Алгоритм C5.0 позволяет назначать штрафы для разных типов ошибок, чтобы воспрепятствовать принятию решений, влекущих за собой более дорогостоящие ошибки. Штрафы обозначены в *матрице штрафов*, в которой указано, во сколько раз одни ошибки обходятся дороже других.

Для построения матрицы штрафов сначала нужно определить ее размеры. Поскольку и прогнозируемые, и фактические классы могут принимать два значения — `yes` или `no`, — нужно создать матрицу 2×2 , используя список из двух векторов, каждый из которых имеет два значения. Одновременно мы присвоим название строкам и столбцам матрицы во избежание путаницы в дальнейшем:

```
> matrix_dimensions <- list(c("no", "yes"),  
c("no", "yes"))> names(matrix_dimensions) <-  
c("predicted", "actual")
```

Как показывает исследование нового объекта, размеры были выбраны правильно:

```
> matrix_dimensions$predicted[1]  
"no"      "yes"$actual[1] "no"      "yes"
```

Затем нужно назначить штрафы за различные типы ошибок, указав четыре значения, заполняющих матрицу. Поскольку R заполняет матрицу по столбцам, сверху вниз, то нужно указать значения в следующем порядке.

1. По прогнозу — нет, фактически — нет.
2. По прогнозу — да, фактически — нет.
3. По прогнозу — нет, фактически — да.
4. По прогнозу — да, фактически — да.

Предположим, что невозвращение кредита обходится банку в четыре раза дороже, чем упущенная возможность. Тогда штрафные значения могут быть определены как:

```
> error_cost <- matrix(c(0, 1, 4, 0), nrow =  
2, dimnames =  
matrix_dimensions)
```

В результате будет создана следующая матрица:

```
>
error_cost      actualpredicted      no  yes
no              0      4      yes  1    0
```

Согласно этой матрице, если алгоритм правильно классифицирует `no` и `yes`, штрафы не назначаются. Ложный отрицательный результат имеет штраф 4, а ложный положительный — 1. Чтобы увидеть, как это влияет на классификацию, применим ее к нашему дереву решений, используя параметр `costs` функции `C5.0()`. В остальном мы будем использовать те же операции, что и раньше:

```
> credit_cost <- C5.0(credit_train[-17],
credit_train$default,
costs
= error_cost)> credit_cost_pred <-
predict(credit_cost, credit_test)>
CrossTable(credit_test$default,
credit_cost_pred,
prop.chisq = FALSE,
prop.c = FALSE, prop.r = FALSE,
dnn =
c('actual default', 'predicted default'))
```

В результате получим следующую матрицу несоответствий:

actual default	predicted default		Row Total
	no	yes	
no	37 0.370	30 0.300	67
yes	7 0.070	26 0.260	33
Column Total	44	56	100

По сравнению с усиленной моделью, эта версия в целом допускает больше ошибок: 37 против 18 % в варианте с усилением. Однако типы ошибок очень различаются. Если предыдущие модели правильно классифицировали только 42 и 61 % невозвращений кредитов, то в этой модели $26 / 33 = 79$ % фактических невозвращений были спрогнозированы правильно. Этот компромисс, приводящий к уменьшению ложных отрицательных результатов за счет увеличения ложных положительных результатов, может быть приемлемым, если точно оценить стоимость ошибок.

Правила классификации

Правила классификации — это знания, представленные в форме логических операторов «если... то», позволяющие присвоить класс немаркированным примерам. Правила классификации описываются в терминах *предпосылок и следствий*, которые формируют утверждение, гласящее:

«Если случится это, то случится и то». Предпосылки содержат определенные сочетания значений признаков, в то время как следствия указывают на значение класса, которое присваивается, если выполняются условия правила. Простое правило может гласить: «Если жесткий диск издает щелчки, значит, вот-вот случится сбой».

Алгоритмы обучения на основе правил часто применяются для подобных задач. Как и деревья решений, они могут использоваться в приложениях, где генерируются знания для будущих действий, в том числе таких, как:

- выявление условий, которые приводят к механическим сбоям;
- описание ключевых характеристик групп людей для разбиения клиентов на группы;
- поиск условий, которые предшествуют значительному падению или росту цен акций на фондовом рынке.

У алгоритмов обучения на основе правил есть ряд особенностей, резко отличающих их от деревьев решений. В отличие от дерева, которому необходимо пройти через последовательность ветвящихся решений, правила — это предложения, которые можно воспринимать как независимые утверждения фактов. Кроме того, по причинам, которые мы обсудим далее, результаты алгоритмов обучения на основе правил могут оказаться более простыми, непосредственными и понятными, чем дерево решений, построенное на тех же данных.



Возможно, вы уже поняли, что ветви деревьев решений — это почти то же самое, что операторы «если... то» в алгоритмах обучения на основе правил, и на самом деле правила могут быть сгенерированы из деревьев. Так зачем возиться с отдельной группой алгоритмов обучения на основе правил? Читайте дальше — и узнаете, чем различаются эти два подхода. Алгоритмы обучения на основе правил обычно применяются для решения задач, в которых большинство или все признаки являются номинальными. Эти алгоритмы хорошо выявляют редкие события, даже если такие события происходят лишь для очень специфического сочетания значений признаков.

Отделяй и властвуй

Алгоритмы обучения на основе правил классификации используют подход, известный как «отделяй и властвуй» (separate and conquer). Этот процесс основан на таком правиле: охватывается некое подмножество примеров из тренировочных данных, а затем происходит отделение этого подмножества от остальных данных. По мере добавления новых правил полученные

подмножества данных отделяются до тех пор, пока не будет охвачен весь набор данных и примеров больше не останется. Метод «отделяй и властвуй» во многом похож на описанный ранее подход «разделяй и властвуй», однако имеет нюансы, которые скоро станут известны.

Один из способов представить процесс обучения правилам методом «отделяй и властвуй» — представить последовательную детализацию данных, создавая все более специфические правила для определения значений классов. Предположим, что вам поручили сформулировать правила, позволяющие определить, является ли животное млекопитающим.

Вы можете представить множество всех животных в виде большого пространства, как показано на рис. 5.7.



Рис. 5.7. Алгоритм обучения на основе правил позволяет разделить всех животных на млекопитающих и не млекопитающих

Алгоритм обучения на основе правил начинается с использования признаков, доступных для поиска однородных групп. Например, используя признак, указывающий на то, передвигается ли данный вид по суше, морю или воздуху, можно построить первое правило, предполагающее, что все наземные животные являются млекопитающими (рис. 5.8).

Увидели ли вы какие-либо несоответствия в этом правиле? Если вы любитель животных, то могли заметить, что лягушки не млекопитающие, а земноводные. Данное правило должно быть более конкретным. Попробуем углубиться в детали и предположим, что млекопитающие должны не только передвигаться по земле, но и иметь хвост (рис. 5.9).

Как видно на рис. 5.9, новое правило приводит к подмножеству животных, все из которых являются млекопитающими. Таким образом, подмножество млекопитающих может быть отделено от других данных, и лягушки в этой эволюции делают шаг назад, без обид!



Рис. 5.8. Потенциальное правило рассматривает животных, которые передвигаются по суше, как млекопитающих



Рис. 5.9. Более конкретное правило предполагает, что животные, которые передвигаются по земле и имеют хвосты, являются млекопитающими

Еще одно правило может быть сформулировано для выделения летучих мышей — единственного оставшегося млекопитающего. Потенциальной особенностью, отличающей летучих мышей от остальных животных, является наличие шерсти. Используя правило, построенное на основе этого признака, мы правильно идентифицировали всех животных, что видно на рис. 5.10.



Рис. 5.10. Правило, гласящее, что животные, имеющие шерсть, являются млекопитающими, прекрасно классифицирует оставшихся животных

Поскольку все обучающие экземпляры были классифицированы, процесс обучения на основе правил на этом этапе останавливается. Из этого следует три вывода:

- животные, которые передвигаются по земле и имеют хвосты, являются млекопитающими;
- если у животного нет шерсти, то это не млекопитающее;
- в противном случае животное является млекопитающим.

В предыдущем примере показано, как правила постепенно охватывают все более обширные сегменты данных, чтобы в итоге классифицировать все экземпляры. Поскольку правила, по-видимому, охватывают части данных, алгоритмы, построенные по схеме «отделяй и властвуй», также известны как *алгоритмы покрытия*, а полученные в результате правила — как *правила покрытия*. Из следующего раздела вы узнаете, как правила покрытия применяются на практике, на примере простого алгоритма обучения на основе правил. Затем рассмотрим более сложный алгоритм обучения на основе правил, оба алгоритма применим к реальной задаче.

Алгоритм 1R

Предположим, что в телеигре за занавесом спрятано животное. Вас просят угадать, млекопитающее это или нет. Если вы угадаете, то выиграете большой денежный приз. Вам не дают никаких подсказок о характеристиках животного, но вы знаете, что очень немногие из существующих в мире животных — млекопитающие. Следовательно, вы отвечаете «не млекопитающее». Что вы думаете о своих шансах?

Подобный ответ, конечно, дает вам максимальные шансы на победу, так как это наиболее вероятный результат, если исходить из предположения, что животное выбрано случайно. Такое телешоу, конечно, нелепо, однако оно демонстрирует простейший классификатор *ZeroR* — алгоритм обучения

на основе правил, который не учитывает никаких признаков и буквально не изучает правила (отсюда и название). Для каждого немаркированного примера, независимо от значений его признаков, этот алгоритм прогнозирует наиболее распространенный класс. Практической пользы от этого алгоритма очень мало, за исключением того, что он обеспечивает простую базу для сравнения с другими, более сложными алгоритмами обучения на основе правил.

Алгоритм *1R* (он же *One Rule*, или *OneR*) улучшен по сравнению с *ZeroR* в том смысле, что выбирает только одно правило. Этот подход может показаться чрезмерно упрощенным, однако он часто работает лучше, чем можно ожидать. Как показали эмпирические исследования, во многих реальных задачах точность этого алгоритма иногда приближается к точности гораздо более сложных алгоритмов.



Для более подробного изучения удивительной эффективности алгоритма *1R* обратитесь к статье: Holte R.C. Very Simple Classification Rules Perform Well on Most Commonly Used Datasets. *Machine Learning*, 1993. Vol. 11. P. 63–91.

В табл. 5.2 перечислены достоинства и недостатки алгоритма *1R*.

Таблица 5.2

Достоинства	Недостатки
Создает одно простое для понимания, доступное человеку правило. Часто на удивление хорошо работает. Может служить ориентиром для сравнения с более сложными алгоритмами	Использует только одну функцию. Вероятно, чрезмерно упрощен

Этот алгоритм работает очень просто. По каждому признаку *1R* делит данные на группы с одинаковыми значениями этого признака. Затем для каждого полученного сегмента алгоритм прогнозирует класс большинства. Для каждого правила, основанного на конкретном признаке, вычисляется коэффициент ошибки. В качестве единственного правила выбирается правило с наименьшим коэффициентом ошибки.

На рис. 5.11 показано, как этот алгоритм работает в случае с рассмотренными ранее данными о животных.

Животное	Передвигается по	Имеет шерсть	Млекопитающее
Летучие мыши	Воздух	Да	Да
Медведи	Земля	Да	Да
Птицы	Воздух	Нет	Нет
Кошки	Земля	Да	Да
Собаки	Земля	Да	Да
Угри	Вода	Нет	Нет
Слоны	Земля	Нет	Да
Рыбы	Вода	Нет	Нет
Лягушки	Земля	Нет	Нет
Насекомые	Воздух	Нет	Нет
Свиньи	Земля	Нет	Да
Кролики	Земля	Да	Да
Крысы	Земля	Да	Да
Носороги	Земля	Нет	Да
Акулы	Вода	Нет	Нет

Полный набор данных

Передвигается по	Прогноз	Млекопитающее
Воздух	Нет	Да
Воздух	Нет	Нет
Воздух	Нет	Нет
Земля	Да	Да
Земля	Да	Да
Вода	Да	Да
Земля	Да	Да
Земля	Да	Нет
Земля	Да	Да
Земля	Да	Да
Земля	Да	Да
Земля	Да	Да
Вода	Нет	Нет
Вода	Нет	Нет
Вода	Нет	Нет

×

×

Имеет шерсть	Прогноз	Млекопитающее
Нет	Нет	Нет
Нет	Нет	Нет
Нет	Нет	Да
Нет	Нет	Нет
Нет	Нет	Нет
Нет	Нет	Нет
Нет	Нет	Да
Нет	Нет	Да
Нет	Нет	Нет
Да	Да	Да
Да	Да	Да
Да	Да	Да
Да	Да	Да
Да	Да	Да
Да	Да	Да

×

×

×

Правило «Передвигается по»
Коэффициент ошибки = 2/15

Правило «Имеет шерсть»
Коэффициент ошибки = 3/15

Рис. 5.11. Алгоритм 1R выбирает одно правило с наименьшим коэффициентом ошибки

По признаку «передвигается по» набор данных разделился на три группы: «воздух», «земля» и «вода». В соответствии с прогнозом животные в группах «воздух» и «вода» не являются млекопитающими, в то время как животные в группе «земля» — млекопитающие. Это привело к двум ошибкам — в примере с летучими мышами и лягушками.

По признаку «имеет шерсть» животные разделились на две группы: те, кто покрыт шерстью, согласно прогнозу, являются млекопитающими, а те, у кого шерсти нет, — не млекопитающие. Здесь возникли три ошибки: свиньи, слоны и носороги.

Поскольку выбор признака «передвигается по» привело к меньшему количеству ошибок, то алгоритм 1R будет возвращать следующий результат:

- если животное передвигается по воздуху, то это не млекопитающее;
- если животное передвигается по суше, то это млекопитающее;
- если животное передвигается по воде, то это не млекопитающее.

На этом алгоритм останавливается, найдя самое важное правило.

Очевидно, что для некоторых задач такой алгоритм обучения на основе правил может оказаться слишком примитивным. Вряд ли вы захотите, чтобы медицинская диагностическая система делала вывод только по одному симптому или автопилот основывался только на одном факторе, принимая решение о том, останавливать или ускорять автомобиль. Для таких задач необходим более сложный алгоритм обучения на основе правил, речь о котором пойдет в следующем разделе.

Алгоритм RIPPER

Первые алгоритмы обучения на основе правил имели две проблемы. Во-первых, они были известны своей медлительностью, что делало их неэффективными в условиях растущего числа больших наборов данных. Во-вторых, они часто были склонны к неточности при зашумленных данных.

Первый шаг к решению этих проблем был предложен в 1994 году Йоханнесом Фернкранцем (Johannes Furnkranz) и Герхардом Уидмером (Gerhard Widmer). Предложенный ими алгоритм *постепенного сокращения ошибок* (Incremental Reduced Error Pruning, IREP) использует комбинацию методов раннего и позднего сокращения, которые вырабатывают очень сложные правила и сокращают их, прежде чем отделять экземпляры от полного набора данных. Эта стратегия позволила повысить эффективность алгоритмов обучения на основе правил, однако деревья решений часто оказывались эффективнее.



Подробнее об алгоритме IREP читайте в статье: Furnkranz J., Widmer G. Incremental Reduced Error Pruning // Proceedings of the 11th International Conference on Machine Learning, 1994. P. 70–77.

Следующий шаг в развитии алгоритмов обучения на основе правил был сделан в 1995 году, когда Уильям У. Коэн (William W. Cohen) ввел *многократное отсечение с приращением для приведения погрешности* (Repeated Incremental Pruning To Produce Error Reduction, RIPPER), что позволило улучшить IREP и генерировать правила, эффективность которых сравнима

с эффективностью алгоритма построения деревьев решений или превышает ее.



Подробнее о RIPPER читайте в публикации: Cohen W.W. Fast Effective Rule Induction // Proceedings of the 12th International Conference on Machine Learning, 1995. P. 115–123.

Как показано в табл. 5.3, по своим достоинствам и недостаткам метод RIPPER, как правило, сопоставим с деревьями решений. Главное преимущество этого алгоритма состоит в том, что он позволяет получить более простую модель.

Таблица 5.3

Достоинства	Недостатки
Создает простые, понятные человеку правила. Эффективен на больших и зашумленных наборах данных. Как правило, создает более простую модель, чем сопоставимое дерево решений	Иногда приводит к правилам, которые, как кажется, игнорируют здравый смысл или экспертные знания. Неидеален для работы с числовыми данными. Иногда работает не так хорошо, как более сложные модели

Алгоритм RIPPER, разработанный после нескольких итераций алгоритмов обучения на основе правил, представляет собой набор эффективных подходов для обучения. Поскольку данный алгоритм достаточно сложный, мы не будем подробно рассматривать его в этой книге. Однако в целом алгоритм можно представить как процесс, состоящий из трех этапов.

1. Рост.
2. Сокращение.
3. Оптимизация.

На этапе роста используется метод «отделяй и властвуй», чтобы «жадно» добавить к правилу новые условия, пока не будет классифицировано подмножество данных или не будут исчерпаны свойства для разделения. Подобно деревьям решений, для идентификации следующего свойства разделения используется критерий получения информации. Если увеличение специфичности правила больше не уменьшает энтропию, правило сразу сокращается. Шаги 1 и 2 повторяются до тех пор, пока не будет достигнут критерий остановки, после чего весь набор правил оптимизируется с использованием различных эвристик.

Алгоритм RIPPER может создавать гораздо более сложные правила, чем алгоритм 1R, поскольку он способен учитывать несколько признаков. Это означает, что алгоритм RIPPER способен создавать правила с несколькими предпосылками, такими как «если животное летает и имеет шерсть, то это млекопитающее». Это улучшает способность алгоритма к моделированию сложных данных, однако, как и в случае деревьев решений, это означает, что правила могут быстро стать трудными для понимания.



Процесс эволюции алгоритмов обучения на основе правил классификации не остановился на RIPPER. Вскоре появились новые алгоритмы обучения на основе правил. В обзорах литературы приводятся алгоритмы IREP ++, SLIPPER, TRIPPER и многие другие.

Правила, построенные на основе деревьев решений

Правила классификации также могут быть получены непосредственно из деревьев решений. Начиная с концевых узлов и следуя по ветвям обратно до корня, можно получить последовательность решений, которые объединяются в одно правило. На рис. 5.12 показано, как можно построить правила классификации из дерева решений для прогнозирования успеха фильма.

Следуя по пути от корневого узла до каждого из концевых узлов, получим такие правила классификации.

1. Если знаменитостей немного, то фильм будет провальным.
2. Если знаменитостей много, а бюджет большой, то фильм станет лидером проката.
3. Если знаменитостей много, а бюджет небольшой, то фильм получит высокую оценку у критиков.

По причинам, которые будут разъяснены в следующем разделе, главный недостаток использования дерева решений для генерации правил заключается в том, что результирующие правила часто оказываются более сложными, чем те, которые можно получить с помощью алгоритма обучения на основе правил. Стратегия «разделяй и властвуй», используемая в алгоритме построения деревьев решений, смещает результаты иначе, чем в алгоритмах обучения на основе правил. Однако иногда в вычислительном отношении эффективнее генерировать правила на основе деревьев.

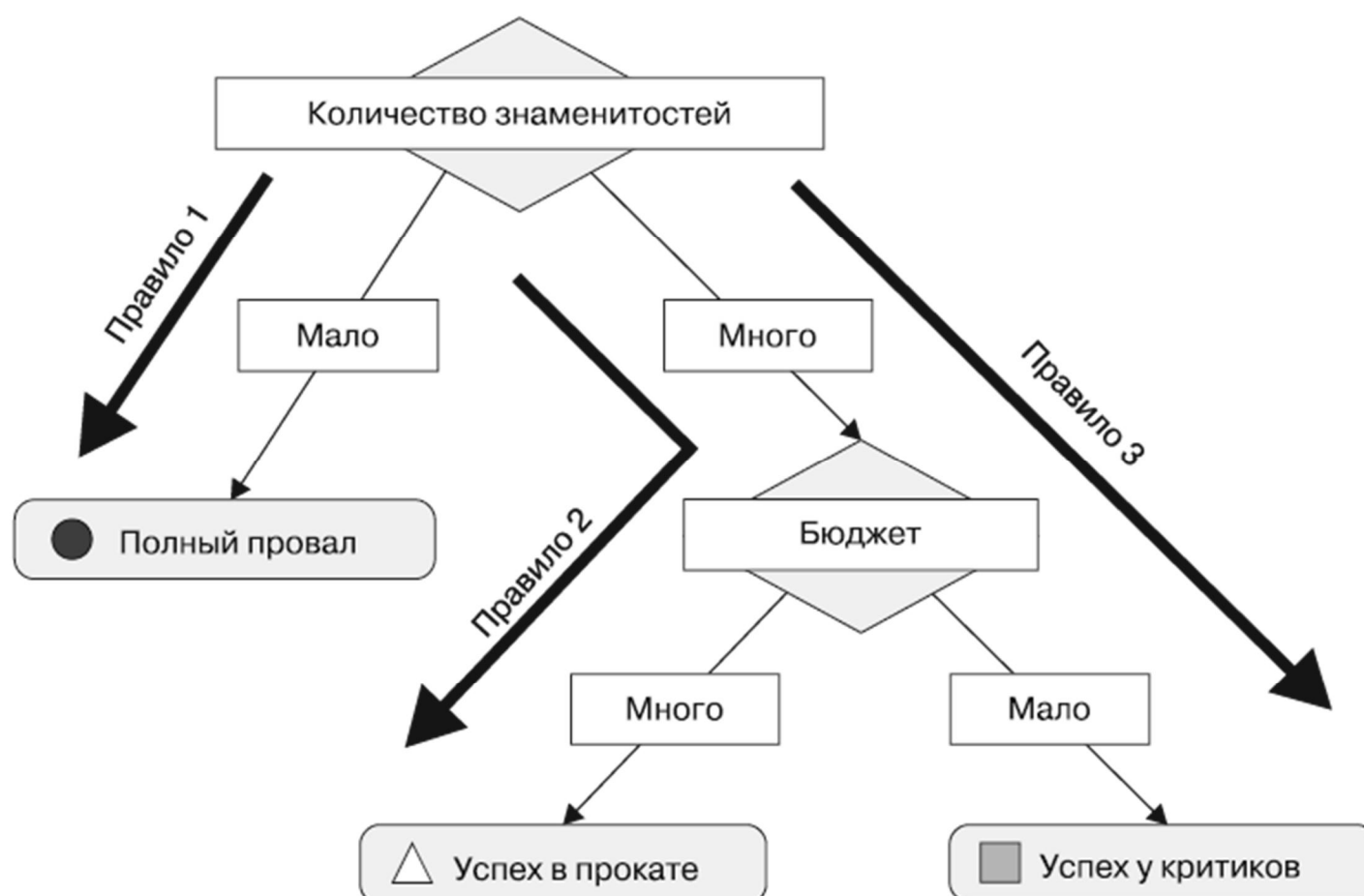


Рис. 5.12. Можно построить правила на основе деревьев решений, следуя по путям от корневого узла до конечных узлов



Если при обучении модели задать значение `rules = TRUE`, то функция `C5.0()` из пакета `C50` будет строить модель с использованием правил классификации.

Когда деревья и правила становятся жадными

Деревья решений и алгоритмы обучения на основе правил известны как *жадные методы обучения*, поскольку обрабатывают данные в порядке поступления. Эвристика «разделяй и властвуй», используемая в деревьях решений, и эвристика «отделяй и властвуй», применяемая в алгоритмах обучения на основе правил, пытаются сгруппировать данные постепенно, сначала выявляя наиболее однородную группу, затем еще более узкую и т.д., пока не будут классифицированы все примеры.

Недостатком жадного подхода является то, что жадные алгоритмы не гарантируют построения оптимального, самого точного или самого короткого списка правил для конкретного набора данных. Принимая сначала самые простые решения, жадный алгоритм обучения может быстро найти одно правило, которое будет точным для какого-то подмножества данных; но при этом он может упустить возможность разработать более детальный набор правил с более высокой общей точностью для всего набора данных. Однако без использования жадного подхода к обучению на основе правил машинное обучение с вычислительной точки зрения было бы применимо разве что для самых маленьких наборов данных (рис. 5.13).

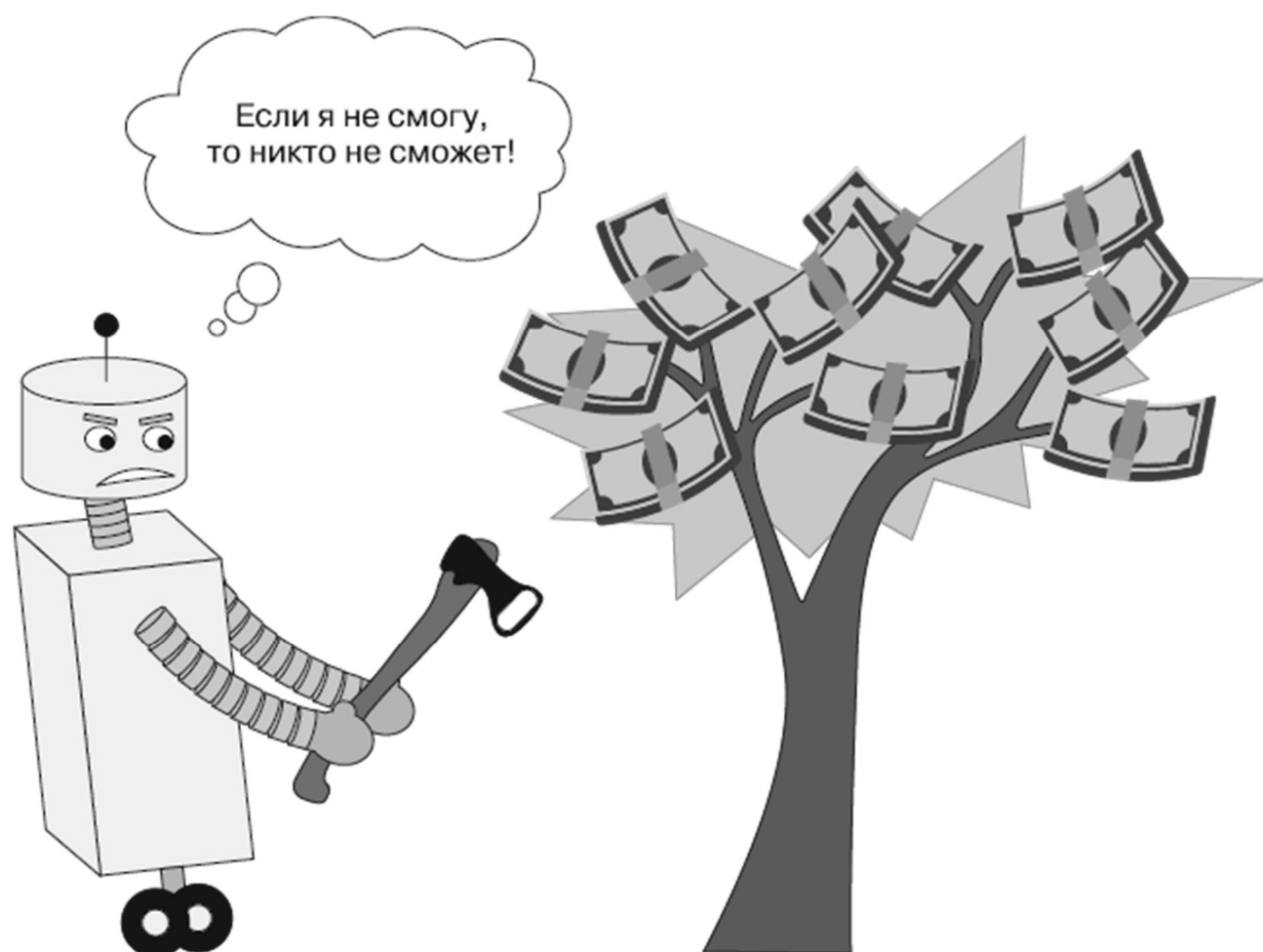


Рис. 5.13. И деревья решений, и изучающие правила классификации являются жадными алгоритмами

Несмотря на то что и деревья, и алгоритмы на основе правил используют жадную эвристику обучения, существуют незаметные различия в том, как строятся эти правила. Возможно, лучший способ их различить — обратить внимание на то, что после разбиения по определенным признакам и формирования подмножеств эти подмножества не могут быть рассмотрены повторно, а только разделены снова. Таким образом, дерево постоянно ограничено историей прошлых решений. И наоборот, если метод «отделяй и властвуй» находит правило, то любые примеры, не охваченные всеми условиями этого правила, могут быть рассмотрены повторно.

Для того чтобы продемонстрировать это различие, обратимся к предыдущему примеру, в котором мы создали алгоритм обучения на основе правил, позволяющий определить, является ли животное млекопитающим. Алгоритм обучения на основе правил определил три правила, которые прекрасно классифицируют животных в данном примере:

- животные, которые передвигаются по земле и имеют хвосты, являются млекопитающими (медведи, кошки, собаки, слоны, свиньи, кролики, крысы и носороги);
- если у животного нет шерсти, то оно не является млекопитающим (птицы, угри, рыбы, лягушки, насекомые и акулы);
- в противном случае животное является млекопитающим (летучие мыши).

Дерево решений, построенное на этих же данных, могло бы вывести четыре следующих правила, что в итоге дало бы такую же идеальную классификацию:

- если животное передвигается по земле и имеет хвост, то это млекопитающее (медведи, кошки, собаки, слоны, свиньи, кролики, крысы и носороги);
- если животное передвигается по земле и не имеет хвоста, то это не млекопитающее (лягушки);
- если животное не передвигается по земле и имеет шерсть, то это млекопитающее (летучие мыши);
- если животное не передвигается по земле и не имеет шерсти, то это не млекопитающее (птицы, насекомые, акулы, рыба и угри).

Разные результаты этих двух подходов объясняются тем, как алгоритм поступит дальше после того, как лягушек отделили по решению «передвигаться по земле». Там, где алгоритм обучения на основе правил допускает повторное рассмотрение лягушек при формировании решения «не имеет шерсти», дерево решений не может модифицировать существующие разделы и, следовательно, вынуждено поместить лягушку в отдельное правило. На рис. 5.14 показан пример обработки результатов.

С одной стороны, поскольку алгоритмы обучения на основе правил позволяют пересматривать примеры, которые уже однажды были рассмотрены, но в итоге не охвачены как часть предыдущих правил, эти алгоритмы часто строят более простой набор правил, чем тот, что строится на основе деревьев решений. С другой стороны, такое повторное использование данных означает, что вычислительные затраты алгоритмов обучения на основе правил могут оказаться выше, чем у деревьев решений.



Рис. 5.14. Различная обработка примера с лягушками эвристиками «разделяй и властвуй» и «отделяй и властвуй»

Пример: распознавание ядовитых грибов по алгоритму обучения на основе правил

Ежегодно много людей получают отравления ядовитыми грибами, некоторые — со смертельным исходом. Поскольку иногда бывает очень сложно отличить ядовитый гриб от съедобного, даже опытные грибники попадают в такие ситуации.

В противоположность вредным растениям, таким как ядовитый дуб или ядовитый плющ, для грибов не существует четких правил идентификации (например, «лист натрое — оставь в покое»), которые позволяли бы определить, является ли гриб ядовитым или съедобным. Дело еще усложняют общепринятые приметы, такие как «все ядовитые грибы ярко окрашены», которые чаще всего являются неверными. Если бы существовали простые, ясные и последовательные правила идентификации ядовитых грибов, они могли бы спасти жизнь грибникам.

Поскольку одним из преимуществ алгоритмов обучения на основе правил является то, что они генерируют простые и понятные правила, эти алгоритмы подходят для решения задачи классификации грибов. Однако правила полезны ровно настолько, насколько они точны.

Шаг 1. Сбор данных

Для того чтобы выделить правила распознавания ядовитых грибов, мы будем использовать набор данных *Mushroom*, составленный Джеффом Шлиммером (Jeff Schlimmer) из Университета Карнеги-Мелона. Необработанный набор этих данных доступен бесплатно в репозитории машинного обучения *UCI Machine Learning Repository* (<http://archive.ics.uci.edu/ml>).

Этот набор данных включает в себя информацию о 8124 грибах из 23 видов пластинчатых грибов, перечисленных в справочнике грибов Северной Америки *Audubon Society Field Guide to North American Mushrooms* (1981). Каждый из видов грибов, приведенных в этом справочнике, классифицируется как «безусловно съедобный», «однозначно ядовитый» или «вероятно, ядовитый, не рекомендуется употреблять в пищу». Для целей нашего набора данных последняя группа была объединена с группой определенно ядовитых грибов, чтобы получить два класса: ядовитые и неядовитые грибы. В словаре данных, доступном на сайте UCI, описаны 22 признака образцов грибов, включая такие характеристики, как форма и цвет шляпки, запах, размер и цвет пластинок, форма ножки и среда обитания.



В этой главе будет использована несколько модифицированная версия данных о грибах. Если вы намерены выполнить этот пример, загрузите файл `mush.csv` и сохраните его в рабочем каталоге `R`.

Шаг 2. Исследование и подготовка данных

Для начала мы воспользуемся функцией `read.csv()`, чтобы импортировать данные для анализа. Поскольку все 22 объекта и целевой класс являются номинальными данными, укажем значение

параметра `stringsAsFactors=TRUE`, чтобы воспользоваться преимуществами автоматического преобразования факторов:

```
> mushrooms <- read.csv("mushrooms.csv",  
stringsAsFactors = TRUE)
```

В результате работы команды `str(mushrooms)` видно, что данные содержат 8124 наблюдения 23 переменных, как и описано в словаре данных. Большая часть результатов `str()` ничем не примечательна, однако стоит отметить одну особенность. Заметили ли вы что-то необычное в переменной `veil_type` в следующей строке?

```
$ veil_type : Factor w/ 1 level "partial": 1 1  
1 1 1 1 ...
```

Если вам показалось странным, что фактор имеет только один уровень, то вы правы. Для этого признака словарь данных приводит два уровня: `partial` и `universal`, однако в наших данных все примеры классифицируются как `partial`. Вполне вероятно, что этот элемент данных был каким-то образом неправильно закодирован. В любом случае, поскольку тип кольца на ножке гриба не меняется в зависимости от выборки, этот признак не предоставляет полезной информации для прогнозирования. Мы удалим эту переменную из анализа, используя следующую команду:

```
> mushrooms$veil_type <- NULL
```

Присваивая вектору `veil_type` значение `NULL`, R удаляет этот признак из фрейма данных `mushrooms`.

Прежде чем двигаться дальше, посмотрим на распределение переменной `type` в нашем наборе данных:

```
>  
table(mushrooms$type) edible poisonous 4208  
3916
```

Примерно 52 % образцов грибов ($N=4208$) съедобны, а 48 % ($N=3916$) — ядовиты.

Для целей нашего эксперимента мы будем считать 8124 образца данных о грибах исчерпывающим набором всех возможных диких грибов. Это важное предположение означает, что нам не нужно выделять часть образцов из тренировочных данных для тестирования. Мы не будем пытаться разработать правила, охватывающие неизвестные виды грибов, лишь попытаемся выделить правила, которые точно описывают все множество известных видов грибов. Поэтому можно построить и протестировать модель на одних и тех же данных.

Шаг 3. Обучение модели на данных

Если бы мы обучали на этих данных гипотетический классификатор ZeroR, что бы он спрогнозировал? Поскольку ZeroR игнорирует все признаки и просто предсказывает наиболее вероятный результат, то на понятном языке его правило гласит: «Все грибы съедобны». Очевидно, что это не очень полезный классификатор: воспользовавшись им, грибник почти в половине случаев заболел бы или умер! Наши правила должны работать намного лучше, чем этот эталонный тест, и давать безопасные советы. В то же время нам нужны простые правила, которые легко запомнить.

Поскольку простые правила все равно могут быть полезны, посмотрим, как простейший алгоритм обучения на основе правил работает с данными о грибах. Для этого обратимся к классификатору 1R — этот алгоритм выделяет единственный признак, который является наиболее прогнозируемым для целевого класса, и использует его для построения правила.

Воспользуемся реализацией 1R, входящей в состав пакета `OneR`, разработанного Хольгером ван Жуаном-Дидрихом (Holger von Joanne-Diedrich) из Ашаффенбургского университета прикладных наук. Это относительно новый пакет, в котором алгоритм 1R реализован на собственном коде R, что обеспечивает скорость и простоту использования. Если у вас еще нет этого пакета, его можно установить с помощью команды `install.packages("OneR")` и загрузить с помощью `library(OneR)`.

Синтаксис алгоритма классификации на основе правил 1R

Использование функции `OneR()` из пакета `OneR`

Построение классификатора:

```
m <- OneR(class ~ predictors, data = mydata)
```

`class` – столбец во фрейме данных `mydata`, который должен быть спрогнозирован;

`predictors` – R-формула, определяющая признаки во фрейме данных `mydata`, которые должны использоваться для прогноза;

`data` – фрейм данных, которому принадлежат `class` и `mydata`.

Функция возвращает объект модели `OneR`, который можно использовать для прогнозирования.

Прогнозирование:

```
p <- predict(m, test)
```

`m` – модель, полученная в результате обучения функцией `OneR()`;

`test` – фрейм данных, содержащий тестовые данные с теми же признаками, что и тренировочные данные, использованные для построения классификатора.

Функция возвращает вектор спрогнозированных значений класса.

Пример:

```
mushroom_classifier <- oneR(type ~ odor + cap_color,  
                             data = mushroom_train)  
mushroom_prediction <- predict(mushroom_classifier,  
                                mushroom_test)
```

Чтобы можно было указать модель для обучения, в функции `OneR()` использован синтаксис R-формулы. В синтаксисе формулы оператор `~` (называемый тильдой) отражает взаимосвязь между целевой переменной и ее прогнозом. Слева от тильды стоит переменная класса, которую нужно спрогнозировать, а справа, разделенные оператором `+`, — признаки, по которым делается прогноз. Если нужно смоделировать отношения между классом `y` и прогностическими признаками `x1` и `x2`, нужно написать формулу `y~x1+x2`. Чтобы включить в модель все переменные, используется символ точки. Например, `y~.` определяет связь между `y` и всеми остальными признаками из набора данных.



Синтаксис R-формулы используется во многих R-функциях. У него есть несколько сильных возможностей для описания отношений между прогностическими переменными. Некоторые из этих возможностей мы рассмотрим в следующих главах. Если же вы стремитесь забежать вперед, почитайте документацию с помощью команды `?formula`.

Используя формулу `type~.` в функции `OneR()`, мы можем учесть в алгоритме все существующие признаки данных о грибах, чтобы спрогнозировать тип гриба:

```
> mushroom_1R <- OneR(type ~ ., data = mushrooms)
```

Чтобы проверить созданные правила, можно ввести имя объекта классификатора:

```
> mushroom_1RCall:OneR.formula(formula = type ~ ., data = mushrooms)Rules:If odor = almond then type = edibleIf odor = anise then type = edibleIf odor = creosote then type = poisonousIf odor = fishy then type = poisonousIf odor = foul then type = poisonousIf odor = musty then type = poisonousIf odor = none then type = edibleIf odor = pungent then type = poisonousIf odor = spicy then type = poisonousAccuracy:8004 of 8124 instances classified correctly (98.52%)
```

Изучив результат команды, мы видим, что для построения правил был выбран признак `odor` — запах. Категории запаха, такие как `almond` (миндаль), `anise` (анис) и т.д., устанавливают правила, по которым алгоритм определяет, является гриб съедобным или ядовитым. Например, если гриб пахнет рыбой (`fishy`), гнилью (`foul`), плесенью (`musty`), имеет едкий (`pungent`) или пряный (`spicy`) запах или же его запах похож на креозот (`creosote`), то такой гриб, скорее всего, ядовит (`poisonous`). И наоборот, согласно прогнозам, съедобные (`edible`) грибы имеют более приятный запах, такой как миндаль или анис, либо же не имеют запаха вовсе. Для практического руководства по сбору грибов эти правила могут быть обобщены простым эмпирическим правилом: «Если гриб неприятно пахнет, то он, скорее всего, ядовит».

Шаг 4. Оценка эффективности модели

В последней строке результатов прогнозирования отмечается, что правила верно прогнозируют съедобность 8004 из 8124 образцов грибов, то есть почти 99 %. Однако алгоритм все еще несовершенен: если ядовитый гриб будет классифицирован как съедобный, есть риск отравления.

Чтобы определить, возможно ли это, рассмотрим матрицу несоответствий прогнозируемых и реальных значений. Для этого нужно сначала сгенерировать прогнозы модели 1R, а затем сравнить их с реальными значениями:

```
> mushroom_1R_pred <- predict(mushroom_1R, mushrooms)> table(actual = mushrooms$type, predicted =
```


mushroom_1R_pred)		predicted		actual	
edible	poisonous	edible	4208	0	p
poisonous	120	3796			

Здесь видно, где наши правила расходятся с действительностью. Согласно прогнозу (данные, представленные в столбцах таблицы), все грибы являются съедобными, в то время как на самом деле (данные в строках таблицы) 4208 грибов являются съедобными, а 3916 — ядовитыми. Изучив таблицу, мы видим, что классификатор 1R не классифицировал какие-либо съедобные грибы как ядовитые, однако он классифицировал 120 ядовитых грибов как съедобные — невероятно опасная ошибка!

С учетом того, что алгоритм обучения использовал только один признак, это достаточно хороший результат; если при сборе грибов избегать неаппетитных запахов, то вы почти гарантированно не попадете в больницу. Тем не менее, когда речь идет о жизни и смерти, слово «почти» не подходит. Не говоря уже о том, что издатель справочника едва ли обрадуется перспективе судебного иска, если его читатели заболеют. Проверим, сможем ли мы добавить еще несколько правил и построить более эффективный классификатор.

Шаг 5. Повышение эффективности модели

Для более сложного алгоритма обучения на основе правил воспользуемся функцией `JRip()` — реализацией алгоритма RIPPER на основе Java. Функция `JRip()` входит в пакет `RWeka`, который, как вы помните, был описан в главе 1, в разделе об установке и загрузке пакетов. Если вы еще не установили этот пакет, вам нужно воспользоваться командой `install.packages("RWeka")`, предварительно установив на компьютере Java в соответствии с инструкциями, специфичными для вашей системы. Выполнив эти операции, загрузите пакет с помощью команды `library(RWeka)`.

Синтаксис алгоритма классификации на основе правил RIPPER

Использование функции JRip() из пакета RWeka

Построение классификатора:

```
m <- JRip(class ~ predictors, data = mydata)
```

class – столбец во фрейме данных mydata, который должен быть спрогнозирован;

predictors – R-формула, определяющая признаки во фрейме данных mydata, которые должны использоваться для прогноза;

data – фрейм данных, которому принадлежат class и mydata.

Функция возвращает объект модели RIPPER, который можно использовать для прогнозирования.

Прогнозирование:

```
p <- predict(m, test)
```

m – модель, полученная в результате обучения функцией JRip();

test – фрейм данных, содержащий тестовые данные с теми же признаками, что и тренировочные данные, использованные для построения классификатора.

Функция возвращает вектор спрогнозированных значений класса.

Пример:

```
mushroom_classifier <- JRip(type ~ odor + cap_color,  
                             data = mushroom_train)  
mushroom_prediction <- predict(mushroom_classifier,  
                                mushroom_test)
```

Как показано в окне синтаксиса, процесс обучения модели JRip() очень похож на обучение модели OneR(). Это одно из приятных преимуществ интерфейса R-формул: синтаксис един для всех алгоритмов, что упрощает сравнение различных моделей.

Теперь обучим алгоритм по правилам JRip(), как это было сделано в случае с OneR(), что позволит алгоритму выделять правила среди всех доступных признаков:

```
> mushroom_JRip <- JRip(type ~ ., data =  
mushrooms)
```

Чтобы проверить правила, введем имя классификатора:

```
> mushroom_JripJRIP rules:=====(odor =  
foul) => type=poisonous (2160.0/0.0)(gill_size =  
narrow) and (gill_color = buff) =>  
type=poisonous (1152.0/0.0)(gill_size = narrow)  
and (odor = pungent) => type=poisonous  
(256.0/0.0)(odor = creosote) => type=poisonous  
(192.0/0.0)(spore_print_color = green) =>  
type=poisonous  
(72.0/0.0)(stalk_surface_below_ring = scaly) and  
(stalk_surface_above_ring = silky) =>
```

```
type=poisonous (68.0/0.0)(habitat = leaves) and  
(gill_attachment = free) and (population =  
clustered) => type=poisonous (16.0/0.0) =>  
type=edible (4208.0/0.0)Number of Rules : 8
```

Классификатор `JRip()` сформировал восемь правил по данным о грибах. Простейший способ прочитать эти правила — представить их как список операторов «если... то», похожих на логику программирования. Первые три правила можно представить следующим образом:

- если у гриба неприятный запах, то этот гриб ядовит;
- если у гриба узкие пластинки, которые пахнут мясом, то этот гриб ядовит;
- если у гриба узкие пластинки с пряным запахом, то этот гриб ядовит.

Наконец, восьмое правило гласит, что любой гриб, который не охвачен предыдущими семью правилами, является съедобным. Следуя принципу нашей логики программирования, это правило можно представить так: «иначе гриб съедобен».

Цифры рядом с каждым правилом указывают на количество экземпляров, охваченных данным правилом, и количество неправильно классифицированных экземпляров. Примечательно, что при использовании этих восьми правил не оказалось ни одного ошибочно классифицированного образца грибов. В результате число экземпляров, охваченных последним правилом, точно равно количеству съедобных грибов в приведенных данных ($N = 4208$).

На рис. 5.15 примерно показано, как применяются эти правила к данным о грибах. Если считать, что большой овал содержит все виды грибов, то алгоритм обучения на основе правил определил признаки или наборы признаков, которые выделяют в этой большой группе однородные сегменты. Прежде всего алгоритм обнаружил большую группу ядовитых грибов, однозначно выделяющуюся неприятным запахом. Затем были обнаружены более мелкие и специфические группы ядовитых грибов. После того как определены правила для каждого из ядовитых грибов, все остальные грибы можно считать съедобными. Благодаря природе каждый сорт грибов оказался настолько уникальным, что классификатор смог обеспечить 100%-ную точность.

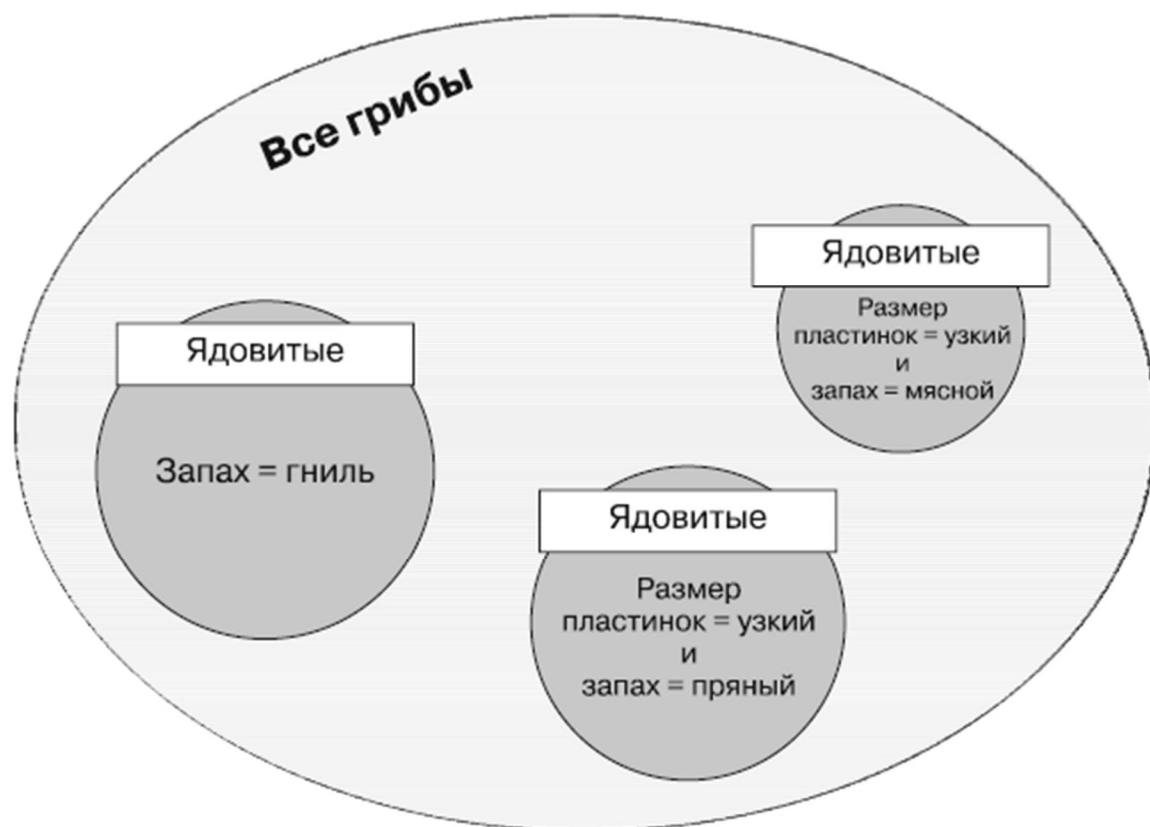


Рис. 5.15. Расширенный алгоритм обучения на основе правил сформулировал правила, полностью охватывающие все виды ядовитых грибов

Резюме

В этой главе были рассмотрены два метода классификации, в которых используются так называемые жадные алгоритмы для разделения данных в соответствии со значениями признаков. Деревья решений используют стратегию «разделяй и властвуй» для создания структур, похожих на блок-схемы, а алгоритмы обучения на основе правил разделяют и группируют данные, чтобы идентифицировать логические правила «если... то». Оба метода позволяют создавать модели, которые можно интерпретировать без статистического обоснования.

Один из популярных настраиваемых алгоритмов построения дерева решений называется C5.0. Мы использовали его для создания дерева, позволяющего спрогнозировать, будет ли возвращен предоставленный кредит. Используя метод усиления и определяя цену ошибок, мы смогли повысить точность алгоритма и избежать выдачи рискованных кредитов, которые могут стоить банку больших денег.

Мы также применили два алгоритма обучения на основе правил, 1R и RIPPER, чтобы сформулировать правила распознавания ядовитых грибов. Алгоритм 1R использовал только один признак, но обеспечил 99%-ную точность идентификации потенциально ядовитых грибов. Однако набор из восьми правил, созданный более сложным алгоритмом RIPPER, правильно определил все съедобные грибы.

В этой главе простыми словами разъяснено, как можно использовать деревья и правила. В следующей главе описываются методы, известные как регрессионные деревья и деревья моделей, которые применяют деревья решений не для классификации, а для числового прогнозирования. В главе 8 увидим, как ассоциативные правила, родственные правилам

классификации, могут использоваться для идентификации групп элементов данных о транзакциях.

Из главы 11 узнаем, как можно повысить эффективность деревьев решений, сгруппировав их в модель, называемую случайным лесом.

6. Прогнозирование числовых данных: регрессионные методы

Математические отношения помогают нам разобраться во многих аспектах повседневной жизни. Например, вес тела зависит от количества потребляемых калорий; доход — от образования и опыта работы; цифры из опросов помогают оценить шансы на переизбрание в президенты.

Если представить такие закономерности в виде чисел, то все станет понятнее. Например, лишние 250 килокалорий, потребляемые ежедневно, могут привести к увеличению веса почти на килограмм в месяц; ежегодно опыт работы растет, что может быть поводом для повышения зарплаты на 1000 долларов в год; а президент, скорее всего, будет переизбран, если экономика стабильна. Очевидно, что эти уравнения не идеально подходят для каждой ситуации, но можно ожидать, что в большинстве случаев они достаточно точны.

Эта глава расширит ваш инструментарий машинного обучения: мы выйдем за рамки описанных ранее методов классификации и рассмотрим методы оценки взаимосвязей между числовыми данными. Исследовав несколько реальных задач числового прогнозирования, вы:

- изучите основные статистические принципы, используемые в регрессии — технологии, которая моделирует размер и прочность числовых взаимосвязей;
- узнаете, как подготовить данные для регрессионного анализа, а также оценить и интерпретировать регрессионную модель;
- познакомитесь с несколькими гибридными методами, известными как регрессионные деревья и деревья моделей, которые адаптируют классификаторы на основе дерева решений для задач числового прогнозирования.

Методы, описанные в этой главе, основаны на большом объеме статистических исследований. Они немного сложнее с математической точки зрения, чем те, которые рассматривались ранее, но не беспокойтесь — даже если вы немного позабыли алгебру, R возьмет на себя самую трудную работу.

Понятие регрессии

Регрессия включает в себя определение взаимосвязи между одной *зависимой числовой переменной* (прогнозируемое значение) и одной или несколькими *независимыми числовыми переменными* (предикторами). Как следует из названия, зависимая переменная зависит от значения независимой

переменной или переменных. Простейшие формы регрессии предполагают, что взаимосвязь между независимой и зависимой переменными является линейной.



Термин «регрессия», применяемый для описания процесса подгонки прямых линий к данным, появился благодаря исследованиям в области генетики сэра Фрэнсиса Гальтона (Sir Francis Galton) в конце XIX века. Гальтон обнаружил, что у низкорослых или, наоборот, очень высоких отцов обычно рождались сыновья, рост которых был ближе к среднему. Он назвал это явление «регрессией к среднему».

Возможно, из уроков алгебры вы помните, что прямую линию можно описать в форме уравнения с угловым коэффициентом вида $y = a + bx$. В этом уравнении y обозначает зависимую переменную, а x — независимую переменную. Угловым коэффициентом указывает, насколько силен наклон прямой по мере увеличения x . Положительные значения соответствуют прямым с наклоном вверх, а отрицательные — с наклоном вниз.

Коэффициент a называют *смещением*, поскольку он определяет точку, в которой прямая пересекает вертикальную ось y . Точка пересечения указывает значение y для $x = 0$. На рис. 6.1 показаны примеры прямых.

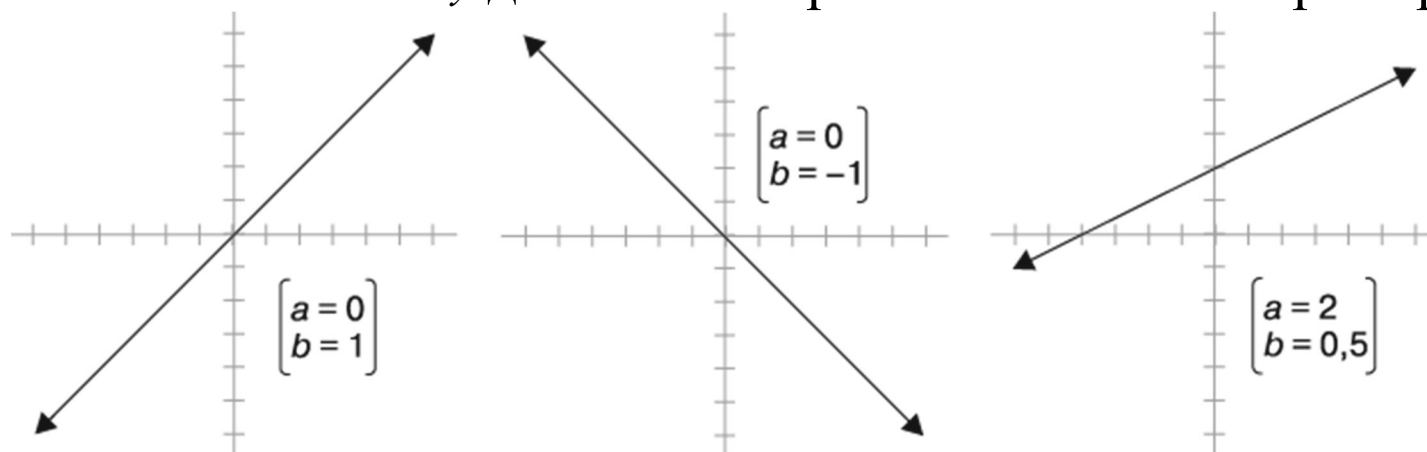


Рис. 6.1. Примеры прямых с разными уклонами и точками пересечения

Данные моделей регрессионных уравнений используют аналогичный формат уравнений с угловым коэффициентом. Задача машины заключается в идентифицировании значений a и b таким образом, чтобы прямая наилучшим образом связывала предоставленные значения x со значениями y . Не всегда удастся построить одну функцию, которая бы идеально связывала эти значения, поэтому у машины также должен быть какой-то способ количественной оценки допустимой погрешности. Вскоре мы поговорим об этом более подробно.

Регрессионный анализ применяется для решения разных задач — почти наверняка это самый распространенный метод машинного обучения. Он может использоваться как для объяснения прошлого, так и для экстраполяции будущего и применим для решения практически любой задачи. Вот некоторые примеры использования этого метода:

- изучение различных измеряемых характеристик для жителей разных стран и отдельных людей в рамках научных исследований в области экономики, социологии, психологии, физики и экологии;
- количественная оценка причинно-следственных связей между событием и его последствиями в таких случаях, как клинические испытания лекарств, технические испытания на безопасность или маркетинговые исследования;
- построение моделей, которые можно использовать для прогнозирования поведения с учетом известных критериев: прогнозирование страховых исков, ущерба от стихийных бедствий, результатов выборов или уровня преступности.

Регрессионные методы также используются для *статистической проверки гипотез* — такая проверка позволяет определить, с какой вероятностью предпосылка является истинной или ложной с учетом наблюдаемых данных. Оценки интенсивности и стабильности отношений, полученные благодаря регрессионной модели, дают информацию, которая может быть использована для того, чтобы определить, являются ли наблюдения случайностью.



Проверка гипотез имеет очень много нюансов и выходит за рамки машинного обучения. Если вас интересует эта тема, советую начать с учебника по введению в статистику: Wolfe D.A., Schneider G. *Intuitive Introductory Statistics*. Springer, 2017.

Регрессионный анализ подразумевает под собой не один алгоритм. Скорее, это общее название для большого числа методов, которые можно адаптировать практически к любой задаче машинного обучения. Если бы вы решили выбрать для изучения только один метод машинного обучения, то регрессия была бы хорошим выбором. Можно посвятить всю жизнь только ей, но все равно всегда будет что изучать.

В этой главе мы сосредоточимся на самых простых моделях *линейной регрессии* — тех, в которых используются прямые линии. Случай, когда существует только одна независимая переменная, называется *простой линейной регрессией*. При наличии двух и более независимых переменных имеет место *множественная линейная регрессия*, или просто *множественная регрессия*. Оба эти метода предполагают наличие единственной зависимой переменной, которая измеряется непрерывно.

Регрессия может также использоваться для других типов зависимых переменных и даже для некоторых задач классификации.

Например, *логистическая регрессия* применяется для моделирования результата события в бинарном виде (0 или 1), а *регрессия Пуассона*, названная в честь французского математика Симона Пуассона (Sime'on Poisson), моделирует результат в виде целых чисел. *Полиномиальная логистическая регрессия* моделирует не бинарный результат и поэтому может использоваться для

классификации. Поскольку во всех регрессионных методах применяются одни и те же статистические принципы, после того как вы освоите линейный случай, изучить остальные варианты не составит труда.



Многие специализированные регрессионные методы относятся к классу *обобщенных линейных моделей* (Generalized Linear Models, GLM). Используя GLM, можно распространить линейные модели на другие шаблоны с помощью *функции связи*, которая определяет более сложные формы зависимости между x и y . Это позволяет применять регрессию практически к любому типу данных.

Начнем с базового случая простой линейной регрессии. Несмотря на название, этот метод не так уж прост при решении сложных задач. Из следующего раздела мы узнаем, как применение модели простой линейной регрессии могло бы предотвратить ужасную катастрофу.

Простая линейная регрессия

Двадцать восьмого января 1986 года семь членов экипажа американского космического челнока «Челленджер» погибли из-за отказа ракетного ускорителя, что привело к его катастрофическому разрушению.

Впоследствии эксперты быстро назвали главную причину трагедии — температура запуска. Резиновые уплотнительные кольца, обеспечивающие герметичность соединений ракеты, никогда не испытывались при температуре ниже 40 °F (4 °C), а погода в день запуска была необычно холодной — ниже нуля.

Оглядываясь назад, можно сказать, что эта авария стала примером того, как важны анализ и визуализация данных. Неизвестно, какая информация была доступна инженерам-ракетостроителям и лицам, принимающим решения о запуске, однако нельзя отрицать, что, будь в их распоряжении более качественные данные, после тщательного анализа катастрофы вполне можно было бы избежать.



Аналитическая часть этого раздела основана на данных, полученных из статьи: Dalal S.R., Fowlkes E.B., Hoadley B. Risk Analysis of the Space Shuttle: Pre-Challenger Prediction of Failure // Journal of the American Statistical Association, 1989. Vol. 84. P. 945–957. Один из вариантов того, как данные могли бы изменить результат, представлен здесь: Tufte E.R., Cheshire C.T. Visual Explanations: Images And Quantities, Evidence And Narrative. Graphics Press, 1997. Противоположная точка зрения изложена в публикации: Representation and misrepresentation: Robison W., Boisjoly R.,

Hoeker D., Young S. Tufte and the Morton Thiokol engineers on the Challenger // Science and Engineering Ethics, 2002. Vol. 8.P. 59–81.

Инженеры-ракетостроители почти наверняка знали, что при низких температурах компоненты становятся более хрупкими и менее пригодными для качественной герметизации, что, вероятно, может стать причиной утечки топлива. Однако, учитывая давление политиков, требовавших продолжать запуск, инженерам нужны были данные, которые бы подтверждали их гипотезу. Здесь им бы очень пригодилась регрессионная модель, демонстрирующая зависимость между температурой и отказами уплотнительных колец, которая могла бы спрогнозировать вероятность отказа с учетом ожидаемой температуры в момент запуска.

Чтобы построить регрессионную модель, ученые могли бы использовать данные о температуре запуска и проблемах с компонентами, записанные во время 23 предыдущих успешных запусков челнока. Разрушение компонентов может указывать на одну из двух проблем. Первая из них — эрозия — возникает, когда избыточное тепло сжигает уплотнительное кольцо. Вторая — продувка — возникает, когда сквозь плохо прилегающее уплотнительное кольцо протекают — «продувают» его — горячие газы. Поскольку у шаттла было в общей сложности шесть первичных уплотнительных колец, за рейс могло происходить до шести аварий. Ракета могла пережить одно или несколько таких разрушений или же быть уничтоженной всего одним, но каждое дополнительное разрушение повышало вероятность катастрофического отказа. На рис. 6.2 показан график разрушений первичных уплотнительных колец, зафиксированных в течение предыдущих 23 запусков, в зависимости от температуры при старте.

Исследуя этот график, можно проследить тенденцию: во время запусков, происходивших при более высоких температурах, было зафиксировано меньше разрушений уплотнительных колец. Кроме того, при самом холодном запуске (53 °F) произошло два таких разрушения — ни при каком другом запуске подобного количества разрушений не наблюдалось. С учетом этой информации тот факт, что «Челленджер» должен был стартовать при температуре ниже более чем на 20 градусов, представляется тревожным. Но насколько сильно следовало беспокоиться? Чтобы ответить на этот вопрос, можно обратиться к простой линейной регрессии.

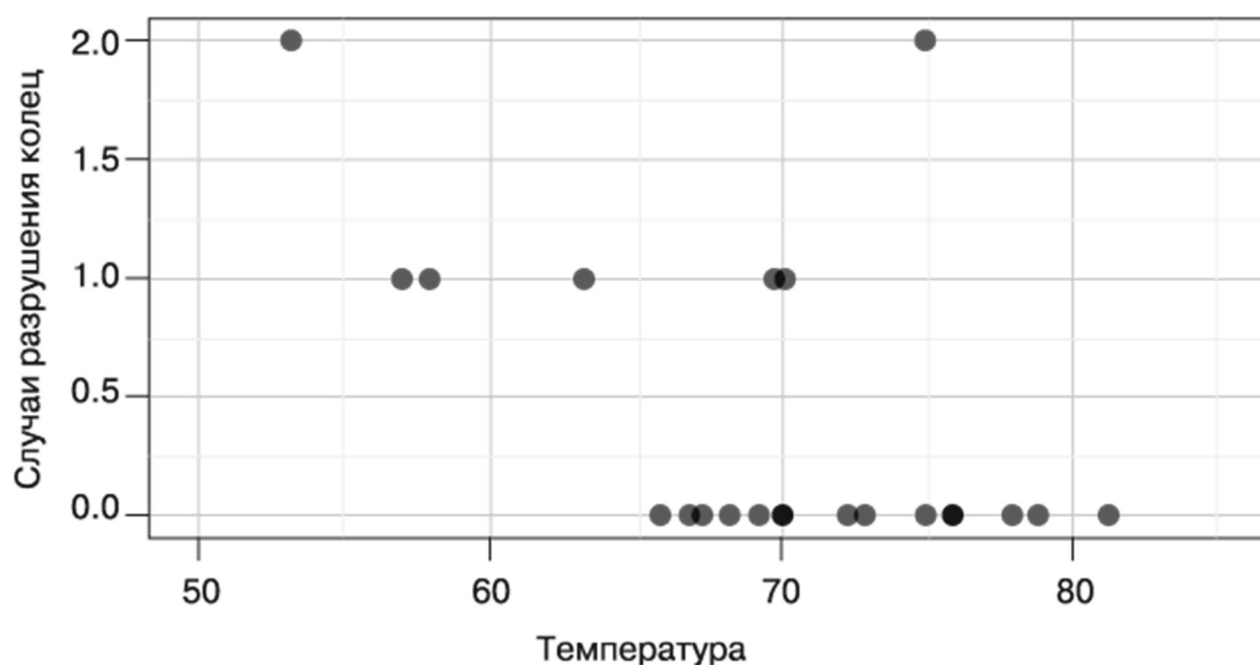


Рис. 6.2. Визуализация повреждений уплотнительных колец космического челнока в зависимости от температуры при старте

Простая линейная регрессионная модель определяет взаимосвязь между зависимой переменной и одной независимой переменной-предиктором с помощью линейной зависимости, описываемой следующим уравнением:
 $y = \alpha + \beta x$.

Если не учитывать греческие буквы, это уравнение практически не отличается от представленного ранее уравнения с угловым коэффициентом. Смещение α (альфа) показывает, в какой точке прямая пересекает ось y , а наклон β (бета) — скорость изменения y при увеличении x . Для данных о запуске шаттла наклон будет сообщать ожидаемое изменение частоты отказов уплотнительных колец при повышении температуры в момент запуска на один градус.



Греческие буквы часто используются в статистике для обозначения переменных, которые являются параметрами статистической функции. Поэтому регрессионный анализ означает оценку значений параметров α и β . Оценки параметров для альфа и бета обычно обозначаются как a и b , хотя иногда некоторые из этих терминов используются как взаимозаменяемые. Предположим, что оценка параметров регрессии в уравнении для данных запуска шаттла известна и составляет $a = 3,70$ и $b = -0,048$. Тогда полное линейное уравнение имеет вид $y = 3,70 - 0,048x$. Пока, не обращая внимания на то, как были получены эти числа, мы можем изобразить прямую, как показано на рис. 6.3.

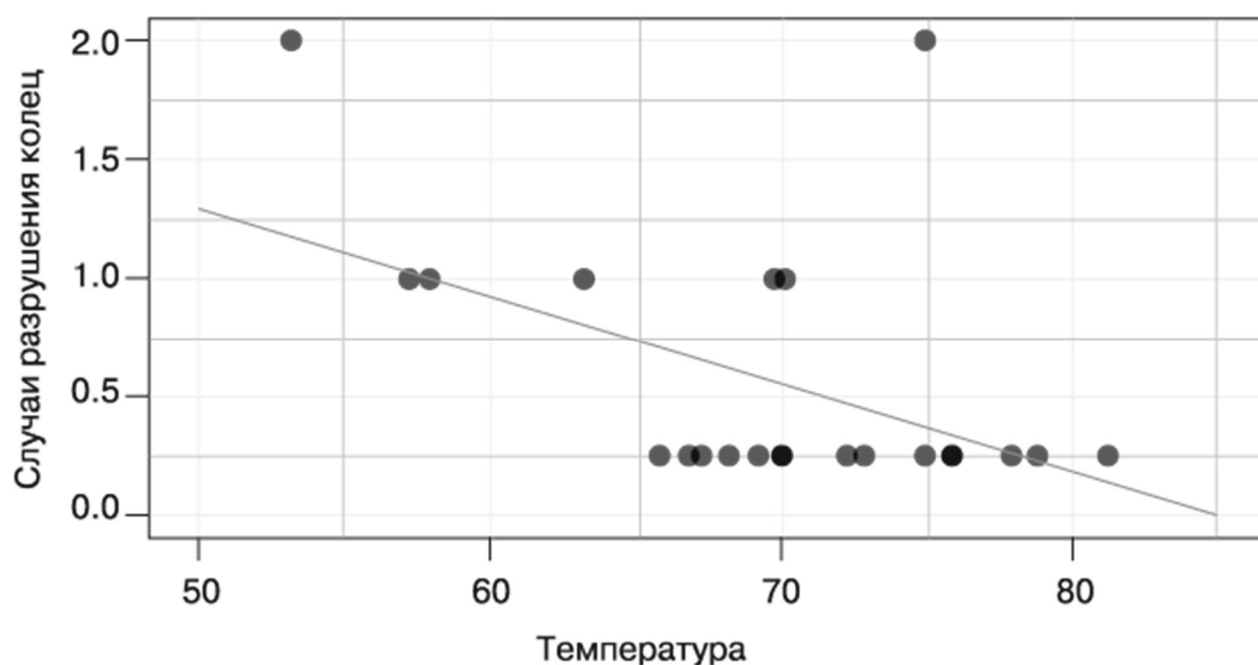


Рис. 6.3. Прямая регрессии, моделирующая зависимость случаев разрушения колец от температуры при запуске челнока

Как видно на рис. 6.3, при 60 °F мы прогнозируем менее одного отказа уплотнительных колец; при 50 °F — ожидаем около 1,3 отказа. Если использовать эту модель для экстраполяции до 31 °F — прогнозируемой температуры при запуске «Челленджера», — то получим примерно $3,70 - 0,048 \times 31 = 2,21$ события отказа уплотнительных колец.

Предполагая, что каждый отказ уплотнительного кольца в равной степени может привести к катастрофической утечке топлива, это означает, что запуск «Челленджера» при 31 °F был почти в три раза более рискованным, чем при обычном запуске при 60 °F, и более чем в восемь раз рискованным, чем запуск при 70 °F.

Обратите внимание, что прямая на графике не проходит точно через все точки данных. Вместо этого она пересекает область данных более или менее равномерно, так что некоторые прогнозы оказываются ниже или выше линии. В следующем разделе мы узнаем, почему была выбрана именно такая линия.

Оценка методом наименьших квадратов

Чтобы вычислить оптимальные оценки α и β , воспользуемся таким методом оценки, как *обычный метод наименьших квадратов* (Ordinary Least Squares, OLS). При OLS-регрессии наклон и смещение выбираются таким образом, чтобы *сумма квадратов ошибок* (Sum of the Squared Errors, SSE) была минимальной. Ошибки, также называемые *невязками*, представляют собой вертикальные расстояния между прогнозируемым значением \hat{y}_i и фактическим значением y_i . Поскольку ошибки могут быть завышенными или заниженными, они могут быть представлены положительными или отрицательными значениями. На рис. 6.4 показаны ошибки для нескольких точек.

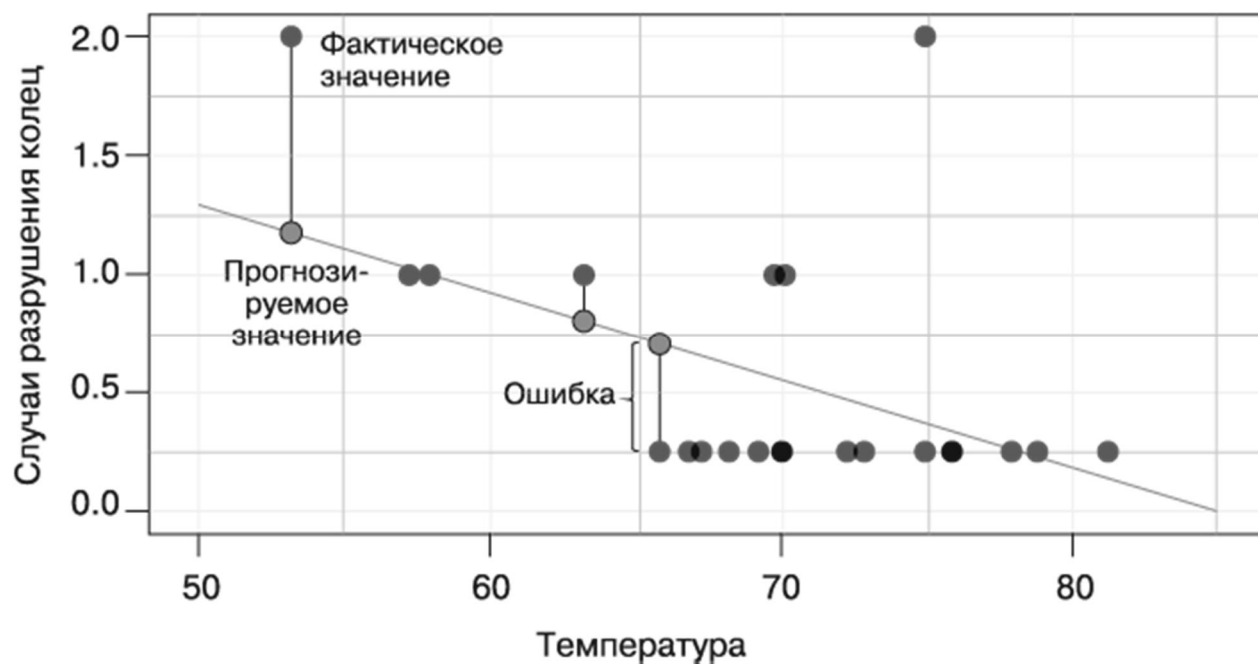


Рис. 6.4. Прогнозируемые значения регрессии отличаются от фактических; разница называется невязкой

В математических терминах цель OLS-регрессии может быть описана как задача минимизации следующего уравнения:

$$\sum (y_i - \hat{y}_i)^2 = \sum e_i^2.$$

Проще говоря, это уравнение определяет e (ошибку) как разницу между фактическим значением y и прогнозируемым значением \hat{y} . Значения ошибок возводятся в квадрат, чтобы исключить отрицательные значения, и суммируются по всем точкам данных.



Знак «крышка» (^) над буквой y является широко распространенным элементом в статистической записи. Этот знак указывает на то, что данная переменная является оценкой истинного значения y . Такое обозначение называется « y с крышкой».

Решение этого уравнения для a зависит от значения b :

$$a = \bar{y} - b\bar{x}.$$



Чтобы понять эти уравнения, вам нужно познакомиться с еще одним обозначением. Горизонтальные полоски над буквами x и y указывают на средние значения x и y . Такие записи называются « x с полоской» и « y с полоской».

Доказательство следующей формулы выходит за рамки данной книги, однако, используя вычисления, можно показать, что значение b , которое приводит к минимальной квадратичной ошибке, вычисляется так:

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}.$$

Если разделить это уравнение на составляющие части, то его можно несколько упростить. Знаменатель должен быть вам знаком; он очень похож на дисперсию x , которая обозначается как $Var(x)$. Как мы знаем из

главы 2, дисперсия вычисляется как среднеквадратичное отклонение от среднего значения x и может быть выражена следующим образом:

$$\text{Var}(x) = \frac{\sum (x_i - \bar{x})^2}{n} .$$

Числитель представляет собой сумму отклонений всех точек данных от среднего значения x , умноженного на отклонение этой точки от среднего значения y . Это похоже на ковариационную функцию для x и y , обозначаемую как $\text{Cov}(x, y)$. Ковариационная формула имеет следующий вид:

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} .$$

Если разделить ковариационную функцию на функцию дисперсии, то n в числителе и знаменателе взаимно сократятся, и можно будет переписать формулу для b следующим образом:

$$b = \frac{\text{Cov}(x, y)}{\text{Var}(x)} .$$

Учитывая это, легко вычислить значение b , используя встроенные R-функции. Применим их к данным о запусках шаттлов, чтобы вычислить регрессионную прямую.



Если вы намерены повторить этот пример, загрузите с сайта Packt Publishing файл `challengeer.csv` и поместите его во фрейм данных с помощью команды `launch <- read.csv("challengeer.csv")`.

Если данные о запусках шаттлов хранятся во фрейме данных с именем `launch`, независимая переменная x называется `temperature`, а зависимая переменная y — `distress_ct`, то для оценки b можно воспользоваться R-функциями `cov()` и `var()`:

```
> b <- cov(launch$temperature,
launch$distress_ct)
/ var(launch$temperature) > b[1] -
0.04753968
```

Затем мы можем оценить a , используя вычисленное значение b и вызвав функцию `mean()`:

```
> a <- mean(launch$distress_ct) - b *
mean(launch$temperature) > a[1] 3.698413
```

Вычислять параметры уравнения регрессии вручную — явно не лучший способ. Поэтому в R встроена функция автоматического подбора регрессионных моделей. В ближайшее время воспользуемся ею. А сейчас важно расширить понимание соответствия регрессионной модели, вначале изучив метод измерения силы линейной зависимости. Кроме того, вы скоро узнаете, как применять множественную линейную регрессию к задачам с несколькими независимыми переменными.

Корреляции

Корреляцией между двумя переменными называется число, которое показывает, насколько близка их взаимосвязь к линейной зависимости. Без дополнительных оговорок корреляцией обычно называют коэффициент корреляции Пирсона, который был введен в XX веке математиком Карлом Пирсоном (Karl Pearson). Корреляция находится в пределах между -1 и $+1$. Максимальное и минимальное значения указывают на идеальную линейную зависимость, в то время как корреляция, близкая к нулю, означает отсутствие линейной зависимости.

Корреляция Пирсона описывается следующей формулой:

$$\rho_{x,y} = \text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} .$$



Здесь введено еще несколько греческих букв: первая из них (которая похожа на строчную букву p) — это «ро», ею обозначают статистическую корреляцию Пирсона. Буквы, похожие на q , повернутую на 45° против часовой стрелки, — это греческие буквы «сигма», ими обозначены стандартные отклонения x и y .

С помощью этой формулы можно вычислить корреляцию между температурой при старте челнока и количеством отказов уплотнительных колец. Напомню, что `cov()` — ковариационная функция, а `sd()` — функция стандартного отклонения. Мы сохраним результат в переменной `r` — этой буквой обычно обозначают предполагаемую корреляцию:

```
> r <- cov(launch$temperature,
launch$distress_ct)
/ (sd(launch$temperature) *
sd(launch$distress_ct))> r[1] -0.5111264
```

Тот же результат можно получить другим способом, с помощью функции корреляции `cor()`:

```
> cor(launch$temperature,
launch$distress_ct)[1] -0.5111264
```

Корреляция между температурой и количеством повреждений уплотнительных колец составляет $-0,51$. Отрицательная корреляция означает, что при повышении температуры количество повреждений уплотнительных колец уменьшается. Для инженеров NASA, изучающих данные об уплотнительных кольцах, это было бы явным показателем того, что запуск при низкой температуре может привести к проблемам. Корреляция также говорит об относительной интенсивности взаимосвязи между температурой и разрушением уплотнительных колец. Поскольку —

0,51 находится посередине между отсутствием корреляции и максимальной отрицательной корреляцией -1 , это означает, что существует умеренно сильная отрицательная линейная зависимость.

Есть ряд эмпирических правил, используемых для интерпретации интенсивности корреляции. Один из методов присваивает статус «слабый» значениям от 0,1 до 0,3, «умеренный» — значениям в диапазоне от 0,3 до 0,5, и «сильный» для значений выше 0,5 (это также относится к аналогичным диапазонам отрицательных корреляций). Однако эти пороговые значения могут быть слишком строгими или слишком слабыми, в зависимости от конкретных целей. Часто корреляция должна интерпретироваться в зависимости от контекста. Для данных, касающихся людей, корреляцию 0,5 можно считать очень высокой; для данных, генерируемых механическими процессами, корреляция 0,5 может быть очень слабой.



Вы, наверное, слышали выражение «корреляция не подразумевает причинно-следственную связь». Дело в том, что корреляция описывает только связь между парой переменных, но могут быть и другие объяснения, которые не подвергались измерениям. Например, может существовать сильная связь между ожидаемой продолжительностью жизни и временем, проведенным за просмотром фильмов, но прежде, чем врачи начнут рекомендовать всем смотреть больше фильмов, необходимо исключить другое объяснение: у молодых людей, которые смотрят много фильмов, вероятность смерти меньше.

Измерение корреляции между двумя переменными позволяет быстро проверить наличие линейной зависимости между независимыми переменными и зависимой переменной. Это станет еще более важным, когда мы будем строить регрессионные модели с большим числом предикторов.

Множественная линейная регрессия

В большинстве случаев при реальном анализе есть несколько независимых переменных. Поэтому вполне вероятно, что во многих задачах числового прогнозирования вы будете использовать множественную линейную регрессию. В табл. 6.1 представлены преимущества и недостатки множественной линейной регрессии.

Таблица 6.1

Преимущества	Недостатки
<p>Самый распространенный в настоящее время подход для моделирования числовых данных.</p> <p>Можно адаптировать для моделирования практически любой задачи.</p> <p>Содержит оценки как размера, так и силы взаимосвязей между признаками и результатами</p>	<p>Делает сильные допущения о данных.</p> <p>Пользователь должен заранее указать форму модели.</p> <p>Не обрабатывает отсутствующие данные.</p> <p>Работает только с числовыми признаками, категориальные данные требуют дополнительной подготовки.</p> <p>Чтобы понять эту модель, нужно разбираться в статистике</p>

Множественную регрессию можно рассматривать как расширение простой линейной регрессии. В обоих случаях цели похожи: найти значения коэффициентов наклона, при которых ошибка прогноза линейного уравнения была бы минимальной. Ключевое отличие множественной регрессии состоит в наличии дополнительных обозначений для независимых переменных.

Модели множественной регрессии имеют вид следующего уравнения. Зависимая переменная y определяется как сумма свободного члена α и каждого из i признаков, представленного переменной x , умноженного на оценочное значение β . Значение ошибки ϵ (обозначается греческой буквой «эпсилон») добавлено для напоминания о том, что ни один прогноз не является идеальным. Ошибка представлена как описанный ранее остаточный член уравнения.

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i + \epsilon.$$

Немного отвлечемся и рассмотрим интерпретацию оценочных параметров регрессии. Как вы заметили, в предыдущем уравнении у каждого признака есть свой коэффициент. Это позволяет любому признаку оказывать свое влияние на значение y . Другими словами, при увеличении признака x на единицу y изменяется на величину β . Следовательно, величина сдвига α является тем ожидаемым значением y , при котором все независимые переменные равны нулю.

Поскольку свободный член уравнения α на самом деле не отличается от остальных параметров регрессии, его иногда обозначают как β_0 (произносится как «бета ноль»), как показано в следующем уравнении:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i + \epsilon.$$

Как и раньше, сдвиг не связан ни с одной из независимых переменных x . Однако по причинам, которые вскоре станут понятны, полезно представить значение β_0 , как если бы оно умножалось на x_0 . Мы будем считать x_0 константой, равной 1.

$$y = \beta_0 x_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i + \epsilon.$$

Для того чтобы оценить параметры регрессии, каждое наблюдаемое значение зависимой переменной y должно быть связано с наблюдаемыми значениями независимых переменных x с помощью представленного выше уравнения регрессии. На рис. 6.5 в графическом виде представлена постановка задачи множественной регрессии.



Рис. 6.5. Цель множественной регрессии — найти значения β , которые связывают значения X и Y при минимальном ε

Строки и столбцы данных на рис. 6.5 могут быть представлены кратко с использованием *матричной записи* — соответствующие переменные выделены полужирным, чтобы подчеркнуть, что в каждой из них содержится несколько значений. Упрощенная таким образом формула выглядит следующим образом:

$$Y = \beta X + \varepsilon.$$

В матричной записи зависимой переменной является вектор Y , содержащий по одной строке для каждого примера. Независимые переменные объединены в матрицу X , число столбцов в которой соответствует числу признаков плюс дополнительный столбец значений «1» для сдвига. В каждом столбце матрицы есть строка для каждого примера. Коэффициент регрессии β и остаточная ошибка ε теперь также являются векторами.

Теперь наша цель состоит в том, чтобы вычислить β — вектор коэффициентов регрессии, который минимизирует сумму квадратов ошибок между прогнозируемыми и фактическими значениями Y . Поиск оптимального решения требует использования матричной алгебры, поэтому вывод формул заслуживает более пристального внимания, но это не предусмотрено в настоящей книге. Однако если вы готовы довериться результатам чужих трудов, то наилучшую оценку вектора β можно рассчитать следующим образом:

$$\hat{\beta} = (X^T X)^{-1} X^T Y.$$

В этом решении используется несколько операций с матрицами: T обозначает *транспонирование* матрицы X , а отрицательный показатель степени указывает на *обратную матрицу*. Используя встроенные матричные операции R , можно реализовать простой обучающий алгоритм

с множественной регрессией. Применим эту формулу к данным о запусках «Челленджера».



Если вы не знакомы с упомянутыми выше матричными операциями, почитайте пару страниц сайта Wolfram MathWorld — <http://mathworld.wolfram.com/Transpose.html> о транспонировании и <http://mathworld.wolfram.com/MatrixInverse.html> — о вычислении обратных матриц. Здесь вы найдете подробное введение и понятное описание, не требующее знаний высшей математики.

Используя следующий код, можно создать простейшую функцию регрессии с именем `reg()`, которая принимает параметры `y` и `x` и возвращает вектор оценочных бета-коэффициентов:

```
reg <- function(y, x) { x <- as.matrix(x) x
<- cbind(Intercept = 1, x) b <- solve(t(x) %*%
x) %*% t(x) %*% y colnames(b) <-
"estimate" print(b)}
```

В созданной функции `reg()` использовано несколько команд R, которые мы раньше не применяли. Во-первых, поскольку мы будем использовать функцию для наборов столбцов из фрейма данных, вызываем функцию `as.matrix()`, которая преобразует фрейм данных в матричную форму. Затем функция `cbind()` связывает с матрицей `x` дополнительный столбец, а по команде `Intercept=1` R присваивает новому столбцу имя `Intercept` и заполняет этот столбец значениями `1`. Затем выполняется последовательность матричных операций над объектами `x` и `y`:

- `solve()` вычисляет обратную матрицу;
- `t()` выполняет транспонирование матрицы;
- `%*%` умножает две матрицы.

Сочетая эти операции согласно формуле, наша функция возвращает вектор `b`, который содержит оценочные параметры для линейной модели, связывающей `x` и `y`. Последние две строки функции присваивают имя вектору `b` и выводят результат на экран.

Применим эту функцию к данным о запусках шаттла. Как показано в следующем коде, набор данных включает в себя три признака и счетчик отказов (`distress_ct`), именно это нас и интересовало:

```
> str(launch)'data.frame':      23 obs. of  4
variables:$ distress_ct      : int    0 1 0 0
0 0 0 0 1 1 ...$ temperature : int    66
```

```

70 69 68 67 72 73 70 57 63 ...$
field_check_pressure : int    50 50 50 50 50 50
100 100 200 ...$ flight_num      : int    1 2
3 4 5 6 7 8 9 10 ...

```

Чтобы убедиться, что функция работает правильно, можно сравнить ее результаты для простой линейной регрессионной модели отказов уплотнительных колец в зависимости от температуры, которая, как мы определили ранее, имеет параметры $a = 3,70$ и $b = -0,048$. Поскольку температура находится во втором столбце данных о запусках, можно выполнить функцию `reg()` с такими параметрами:

```

> reg(y = launch$distress_ct, x =
launch[2])           estimateIntercept    3.6
9841270temperature -0.04753968

```

Эти значения в точности соответствуют нашему предыдущему результату, поэтому используем данную функцию для построения модели множественной регрессии. Применим ее так же, как и раньше, но на этот раз укажем для параметра `x` столбцы со второго по четвертый, чтобы добавить еще два предиктора:

```

> reg(y = launch$distress_ct, x =
launch[2:4])           estimateIntercept
3.527093383temperature
-
0.051385940field_check_pressure  0.001757009flight
t_num          0.014292843

```

Эта модель прогнозирует количество отказов уплотнительных колец, учитывая температуру, давление в полевых условиях и номер запуска. Примечательно, что включение двух дополнительных предикторов не изменило результат, полученный из простой линейной регрессионной модели. Как и прежде, коэффициент для переменной температуры является отрицательным. Это говорит о том, что при увеличении температуры количество ожидаемых отказов уплотнительных колец уменьшается. Величина эффекта также примерно одинакова: на каждый градус повышения температуры при старте ожидается примерно на 0,05 меньше отказов.

Два новых предиктора также вносят вклад в прогнозируемые события отказа. Давление в полевых условиях — это величина давления, приложенного к уплотнительному кольцу во время предпусковых испытаний. Первоначально давление в полевых условиях составляло 50 фунтов на квадратный дюйм, однако при некоторых запусках оно было повышено до 100 и 200 фунтов на квадратный дюйм. Это привело некоторых исследователей к мысли, что именно давление могло стать

причиной разрушения уплотнительного кольца. Коэффициент для данного параметра — положительный, но небольшой, что в некоторой степени служит доказательством этой гипотезы. Номер запуска определяет срок использования челнока.

С каждым полетом ракета изнашивается, детали становятся более хрупкими или подверженными поломке. Небольшая положительная связь между номером рейса и количеством отказов может отражать этот факт.

В целом ретроспективный анализ данных о космическом челноке показывает, что были основания предполагать: с учетом погодных условий запуск «Челленджера» был очень рискованным. Возможно, если бы инженеры тогда применили линейную регрессию, катастрофу можно было бы предотвратить. Но, конечно же, реальная ситуация и связанные с ней политические события были далеко не такими простыми, какими они представляются спустя время.

В этом исследовании лишь поверхностно рассказано о том, что можно делать с помощью линейного регрессионного моделирования. Прделанная работа полезна для точного понимания того, как строятся регрессионные модели, однако при моделировании сложных явлений необходимо учитывать больше условий. Встроенные в R функции регрессии включают дополнительные возможности, требуемые для соответствия этим более сложным моделям, и предоставляют диагностическую информацию, которая помогает в интерпретации модели и оценке ее пригодности. А теперь применим эти функции и расширим знания о регрессии — решим более сложную задачу обучения.

Пример: прогнозирование медицинских расходов с помощью линейной регрессии

Для того чтобы медицинская страховая компания могла зарабатывать деньги, необходимо, чтобы сумма ежегодных взносов превышала расходы на медицинское обслуживание бенефициаров. Следовательно, страховщики вкладывают много времени и денег в разработку моделей, которые точно прогнозируют медицинские расходы застрахованного населения.

Медицинские расходы трудно оценить, поскольку самые дорогостоящие случаи происходят редко и кажутся случайными. Тем не менее некоторые ситуации являются более распространенными для определенных слоев населения. Например, рак легких чаще встречается у курильщиков, чем у некурящих, а от болезней сердца чаще страдают тучные люди.

Целью этого анализа является использование данных о пациентах для прогнозирования средних расходов на медицинское обслуживание для подобных групп населения. Эти оценки могут быть использованы для создания страховых таблиц, согласно которым сумма ежегодных взносов

устанавливается выше или ниже в зависимости от ожидаемых затрат на лечение.

Шаг 1. Сбор данных

Для анализа мы воспользуемся имитационным набором данных, содержащим предполагаемые медицинские расходы для пациентов, проживающих в Соединенных Штатах. Данные были созданы для этой книги с использованием демографической статистики, предоставленной Бюро переписи населения США, и, таким образом, приблизительно соответствуют реальным условиям.



Если вы собираетесь воспроизвести этот пример, загрузите файл `insurance.csv` и сохраните его в рабочей папке R.

В файле `insurance.csv` перечислено 1338 бенефициаров, которые зарегистрированы в программе страхования, с признаками, соответствующими характеристикам пациента, а также общие медицинские расходы, входящие в программу страхования за календарный год. В файле указаны следующие признаки пациентов:

- `age` — целое число, обозначающее возраст основного бенефициара (за исключением лиц старше 64 лет, так как их расходы на медицину обычно покрываются правительством);
- `sex` — пол страхователя: мужской или женский;
- `bmi` — индекс массы тела (ИМТ, или BMI), который позволяет определить, имеет человек недостаточный или избыточный вес. ИМТ вычисляется как вес человека (в килограммах), разделенный на его рост (в метрах) и возведенный в квадрат. Идеальный ИМТ находится в пределах от 18,5 до 24,9;
- `children` — целое число, обозначающее количество детей или иждивенцев, на которых распространяется программа страхования;
- `smoker` — категориальная переменная, принимающая значения «да» или «нет» и указывающая на то, курит ли застрахованное лицо;
- `region` — место жительства получателя страховки в США; разделено на четыре географических региона: северо-восток, юго-восток, юго-запад или северо-запад.

Важно подумать о том, каким образом оплачиваемые медицинские расходы могут зависеть от этих переменных. Например, можно ожидать, что пожилые люди и курильщики сильнее подвержены риску больших медицинских расходов. В отличие от многих других методов машинного обучения при регрессионном анализе зависимость между признаками обычно определяется пользователем, а не распознается автоматически. Некоторые из потенциальных зависимостей будут рассмотрены в следующем разделе.

Шаг 2. Исследование и подготовка данных

Чтобы загрузить данные для анализа, как и раньше, воспользуемся функцией `read.csv()`. Мы можем спокойно использовать параметр `stringsAsFactors=TRUE`, потому что три номинальные переменные целесообразно преобразовать в факторы:

```
> insurance <- read.csv("insurance.csv",  
stringsAsFactors = TRUE)
```

Функция `str()` подтверждает, что данные отформатированы так, как ожидалось:

```
> str(insurance) 'data.frame': 1338 obs. of 7  
variables: $ age : int 19 18 28 33 32 31 46  
37 37 60 ... $ sex : Factor w/ 2 levels  
"female", "male": 1 2 2 2 2 1 ... $ bmi : num  
27.9 33.8 33 22.7 28.9 25.7 33.4 27.7 ... $  
children : int 0 1 3 0 0 0 1 3 2 0 ... $  
smoker : Factor w/ 2 levels "no", "yes": 2 1 1 1  
1 1 1 1 ... $ region : Factor w/ 4 levels  
"northeast", "northwest", ..: ... $ expenses : num  
16885 1726 4449 21984 3867 ...
```

Зависимая переменная в нашей модели — это `expenses`, затраты на медицинское обслуживание, которые покрываются медицинской страховкой в течение года для каждого человека. Перед тем как строить регрессионную модель, полезно проверить нормальность данных. Для линейной регрессии зависимая переменная не обязательно должна иметь нормальное распределение, однако зачастую модель получается лучше, когда это условие выполняется. Посмотрим на сводную статистику:

```
> summary(insurance$expenses) Min. 1st  
Qu. Median Mean 3rd  
Qu. Max. 1122 4740 9382 13270 166  
40 63770
```

Как видим, среднее значение больше медианы. Это означает, что распределение расходов на страхование имеет сдвиг вправо. Наглядно это изображено на рис. 6.6.

```
> hist(insurance$expenses)
```

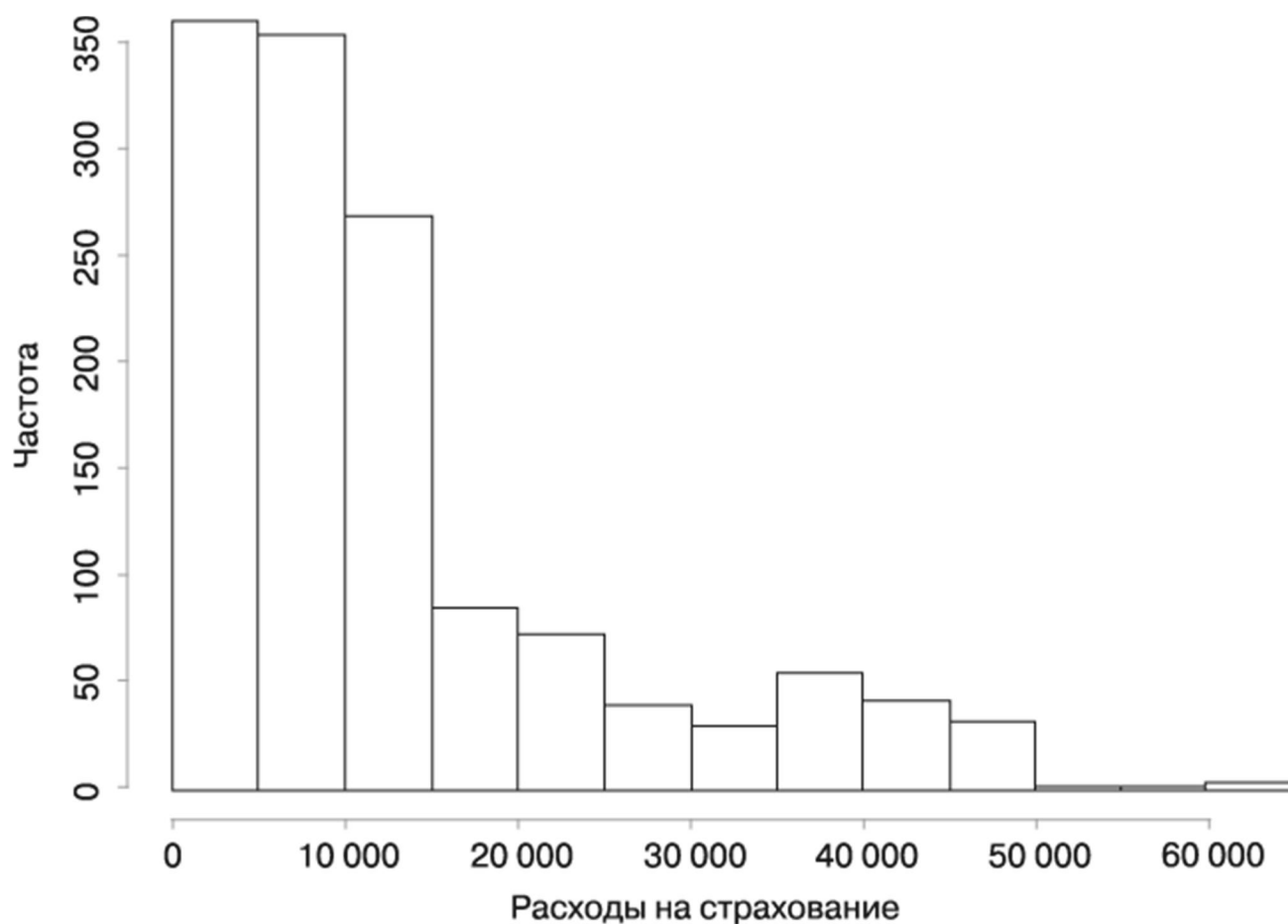


Рис. 6.6. Распределение годовых расходов на страхование

На рис. 6.6 видно, что распределение сдвинуто вправо. Это говорит о том, что у большинства людей в наших данных ежегодные медицинские расходы составляют от 0 до 15 000 долларов, однако хвост распределения простирается далеко за пределы этих пиковых значений. Такое распределение неидеально для линейной регрессии, однако то, что мы заранее знаем об этом недостатке, может помочь нам впоследствии разработать более подходящую модель.

Прежде чем решать эту проблему, обратим внимание на еще одну. Для регрессионных моделей требуется, чтобы все признаки были числовыми, однако в нашем фрейме данных есть три факторных элемента. В частности, переменная `sex` принимает значения `male` и `female`, а переменная `smoker` имеет категории `yes` и `no`. Из результатов `summary()` мы знаем, что признак `region` имеет четыре уровня, и нам стоит присмотреться внимательнее к тому, как они распределяются:

```
>
```

```
table(insurance$region)
northeast    northwest
southeast    southwest    324    325
    364    325
```

Как видим, данные распределяются почти поровну между четырьмя географическими регионами. Вскоре мы узнаем, как R-функция линейной регрессии обрабатывает эти факторные переменные.

Изучение зависимостей между признаками: матрица корреляции
 Прежде чем подстраивать регрессионную модель под данные, иногда полезно определить, каким образом независимые переменные связаны с

зависимой переменной и друг с другом. Быстрый ответ на этот вопрос дает *матрица корреляции*. Она показывает корреляцию для каждой пары переменных из заданного набора.

Для того чтобы создать матрицу корреляции для четырех числовых переменных из фрейма данных о страховании, воспользуемся командой `cor()`:

```
> cor(insurance[c("age", "bmi", "children",
"expenses")])
      age      bmi      children      expenses
children  expenses age      bmi      children      expenses
0.04246900 0.29900819 1.0000000 0.10934101
0.04246900 0.01264471 0.19857626 1.00000000 0.06799823
0.2990082 0.19857626 0.06799823 1.00000000
```

На пересечении каждой строки и столбца находится корреляция для пары переменных, соответствующих этой строке и столбцу. Значение на диагонали всегда равно `1.0000000`, поскольку корреляция переменной к самой себе всегда идеальна. Значения, расположенные симметрично относительно диагонали, идентичны, так как эти корреляции одинаковы. Другими словами, `cor(x, y)` равно `cor(y, x)`.

Ни одна из корреляций в данной матрице не является очень сильной, однако есть несколько заметных зависимостей. В частности, по-видимому, имеется слабая положительная корреляция между `age` и `bmi`, а это означает, что с возрастом масса тела увеличивается. Есть также положительная корреляция между признаками `age` и `expenses`, `bmi` и `expenses`, а также `children` и `expenses`. Эти зависимости означают, что с возрастом, увеличением массы тела и рождением детей ожидаемая стоимость страхования возрастает. Попытаемся выявить эти зависимости более четко, когда построим окончательную модель регрессии.

Визуализация взаимосвязей между признаками: матрица рассеяния
Полезно визуализировать отношения между числовыми объектами с помощью диаграмм рассеяния. Мы могли бы построить диаграмму рассеяния для каждой возможной пары переменных, однако при большом количестве признаков это утомительно.

Вместо этого можно построить *матрицу рассеяния* (сокращенно SPLOM, от ScatterPLOt Matrix), которая представляет собой набор диаграмм рассеяния, представленных в виде сетки. Матрица рассеяния применяется для обнаружения закономерностей среди трех и более переменных. Матрица рассеяния не является настоящей многомерной визуализацией,

поскольку одновременно рассматриваются только два признака. Однако она дает общее представление о взаимосвязях между данными.

Для построения матрицы рассеяния для четырех числовых признаков — `age`, `bmi`, `children` и `expenses` — воспользуемся графическими возможностями R. Функция `pair()`, входящая в состав стандартного пакета R, предоставляет базовые функциональные возможности для построения матриц рассеяния. Для того чтобы вызвать эту функцию, просто предоставьте ей фрейм данных для построения диаграмм. Мы выбрали из фрейма данных о страховании только четыре интересующие нас числовые переменные:

```
> pairs(insurance[c("age", "bmi", "children",  
"expenses")])
```

На рис. 6.7 показано, что получится в результате.

В матрице рассеяния на пересечении каждой строки и столбца содержится диаграмма рассеяния для пары переменных, соответствующих данным строке и столбцу. Диаграммы, расположенные выше и ниже диагонали, представляют собой транспозиции, поскольку для них оси X и Y меняются местами.

Заметили ли вы какие-либо закономерности на этих диаграммах? Некоторые из них выглядят как случайные облака точек, однако на других, похоже, заметны некоторые тенденции. Взаимосвязь между переменными `age` и `expenses` представлена несколькими относительно прямыми линиями, в то время как в зависимости `expenses` от `bmi` просматриваются две группы точек. На остальных диаграммах трудно обнаружить какие-либо тенденции.

Если добавить к диаграммам больше информации, то их можно сделать еще полезнее. С помощью функции `pair.panels()` из пакета `psych` можно построить расширенную матрицу рассеяния. Если у вас не установлен этот пакет,

введите `install.packages("psych")` и загрузите его с помощью команды `library(psych)`. Затем можно построить матрицу диаграммы рассеяния, как мы это уже делали ранее:

```
> pairs.panels(insurance[c("age", "bmi",  
"children", "expenses")])
```

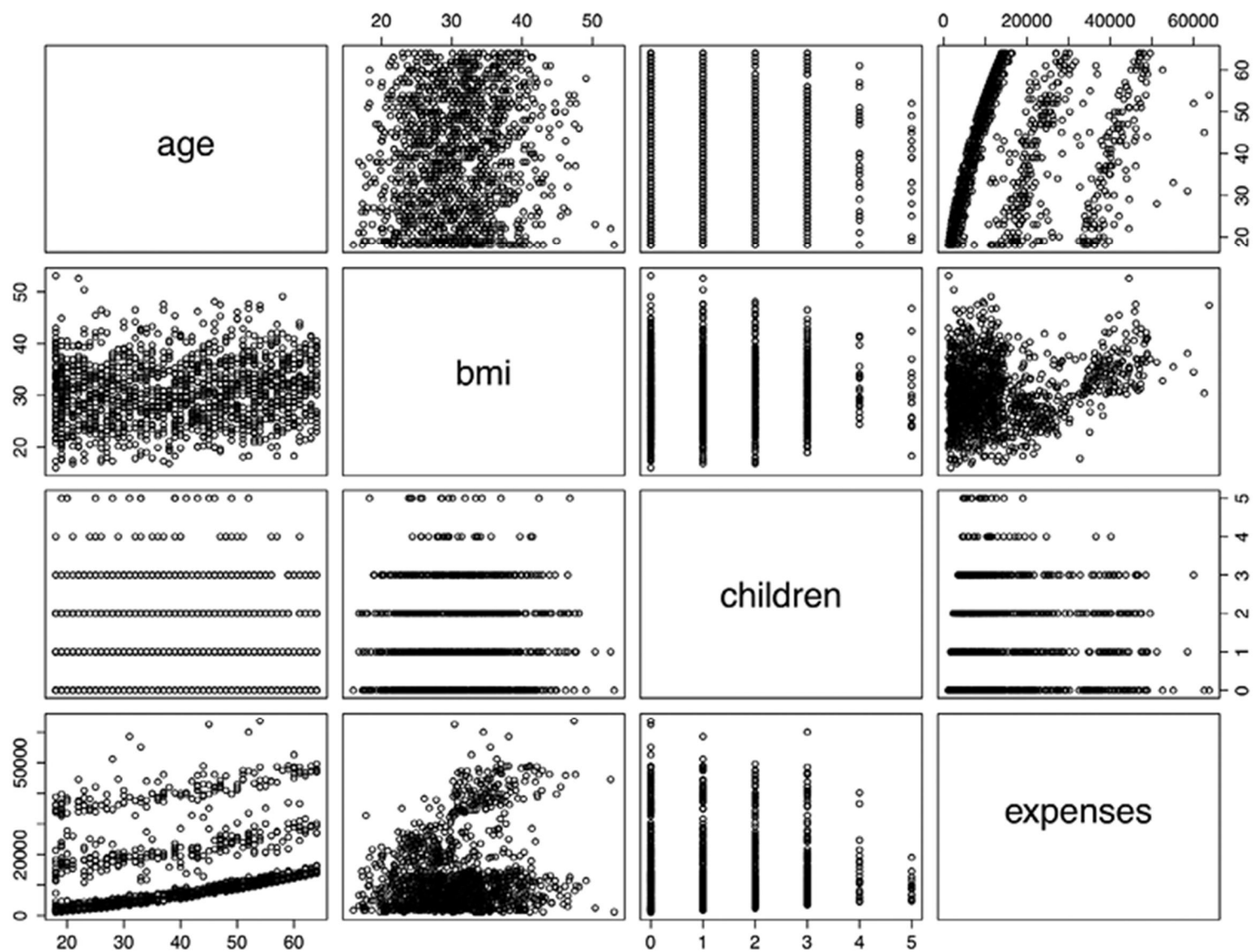


Рис. 6.7. Матрица рассеяния для числовых признаков из набора данных о страховании

Это даст нам немного более информативную матрицу рассеяния, как показано на рис. 6.8.

В результатах `pair.panels()` диаграммы рассеяния представленные выше диагонали заменяются матрицами корреляции. Теперь на диагонали содержатся гистограммы, отображающие распределение значений для каждого признака. Наконец, диаграммы рассеяния, расположенные ниже диагонали, представлены с дополнительной визуальной информацией.

Объект овальной формы, присутствующий на каждой диаграмме рассеяния, называется *эллипсом корреляции*. Он визуализирует силу корреляции. Чем сильнее растянут эллипс, тем сильнее корреляция. Почти идеальный круг, как у зависимости между `bmi` и `children`, указывает на очень слабую корреляцию (в данном случае 0,01).

Эллипс для зависимости между `age` и `expenses` вытянут больше, что говорит о сильной корреляции (0,30). Точка в центре эллипса отражает средние значения переменных по осям x и y .

Кривые, изображенные на диаграммах рассеяния, называются *LOESS-кривыми*. Они указывают на общую зависимость между переменными осей x и y . Это лучше всего объяснить на примере. Кривая для переменных `age` и `children` имеет вид перевернутой буквы U, пик которой соответствует среднему возрасту. Это означает, что самые старые и самые молодые люди в данной выборке имеют меньше детей, которые учитываются в программе страхования, чем люди среднего возраста. Поскольку эта тенденция является нелинейной, то данный вывод не мог

быть сделан на основании одних только корреляций. А вот LOESS-кривая для `age` и `bmi`, наоборот, представляет собой прямую с незначительным наклоном, что означает, что с возрастом масса тела увеличивается, — впрочем, эту зависимость мы уже вывели на основании матрицы корреляции.

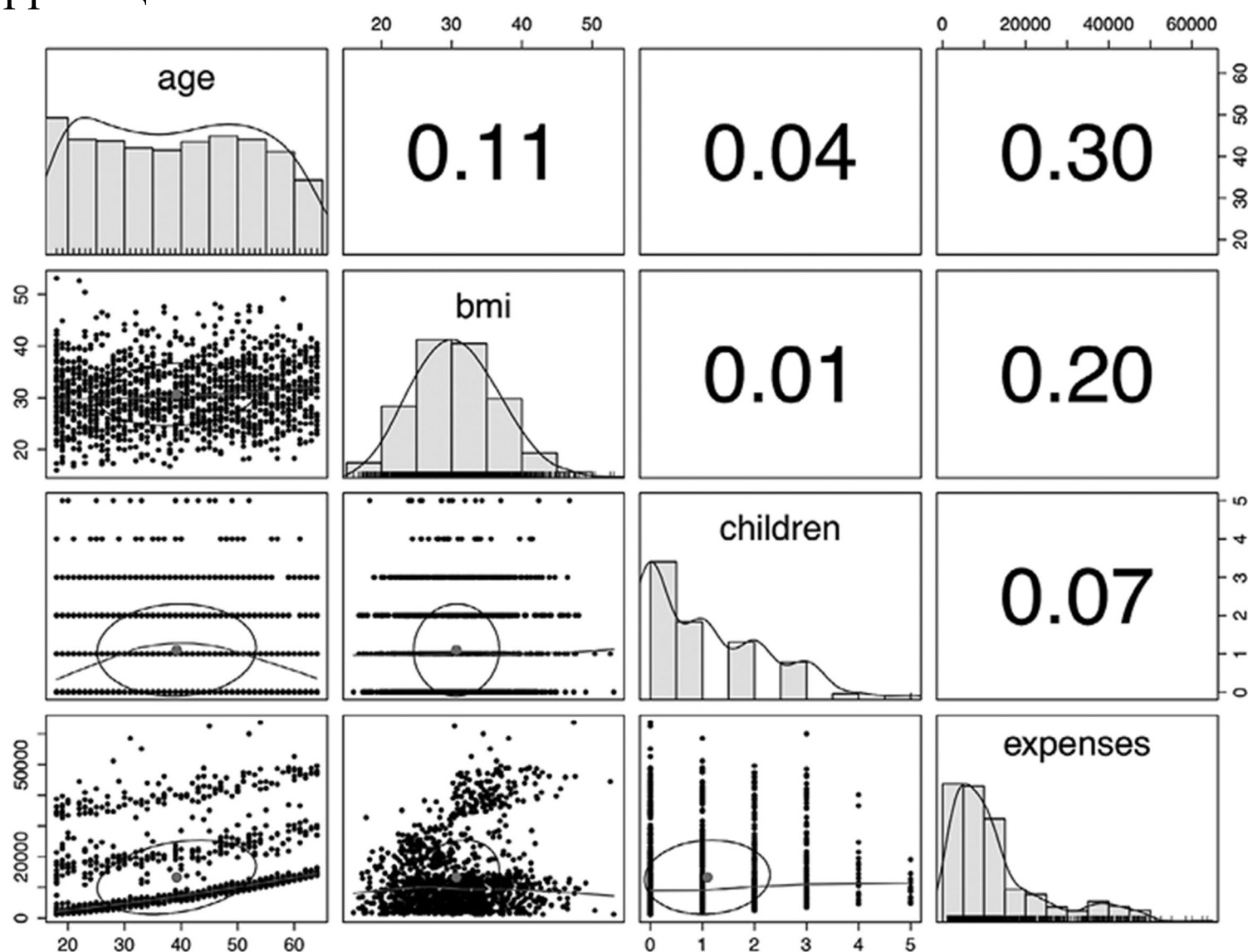


Рис. 6.8. Функция `pair.panels()` добавляет детали к матрице рассеяния

Шаг 3. Обучение модели на данных

Чтобы подобрать модель линейной регрессии для данных с помощью R, можно использовать функцию `lm()`. Она входит в состав пакета `stats`, который входит в стандартный пакет R и загружается по умолчанию. Синтаксис `lm()` выглядит следующим образом.

Синтаксис моделирования множественной регрессии

Использование функции `lm()` из пакета `stats`

Построение модели:

```
m <- lm(dv ~ iv, data = mydata)
```

`dv` – зависимая переменная из фрейма данных `mydata`, которую нужно смоделировать;

`iv` – R-формула, определяющая независимые переменные во фрейме данных `mydata`, которые следует использовать в модели;

`data` – фрейм данных, которому принадлежат переменные `dv` и `iv`.

Функция возвращает объект, который представляет собой регрессионную модель и может быть использован для прогнозирования. Взаимодействия между независимыми переменными обозначаются символом `*`.

Прогнозирование:

```
p <- predict(m, test)
```

`m` – модель, обученная с помощью функции `lm()`;

`test` – фрейм данных, содержащий тестовые данные с теми же признаками, что и тренировочные данные, использованные при построении модели.

Функция возвращает вектор прогнозируемых значений.

Пример:

```
ins_model <- lm(charges ~ age + sex + smoker,  
               data = insurance)  
ins_pred <- predict(ins_model, insurance_test)
```

Следующая команда позволяет построить модель линейной регрессии, которая устанавливает зависимость суммарных медицинских расходов от шести независимых переменных. В синтаксисе R-формулы для описания модели использован символ «тильда» (`~`); слева от тильды стоит зависимая переменная `expenses`, а справа — независимые переменные, разделенные знаками `+`. Указывать величину сдвига для регрессионной модели нет необходимости, так как он учитывается по умолчанию:

```
> ins_model <- lm(expenses ~ age + children +  
bmi + sex + smoker + region, data = insurance)
```

Поскольку для указания всех признаков (кроме тех, которые уже указаны в формуле) может использоваться символ точки (`.`), следующая команда эквивалентна предыдущей:

```
> ins_model <- lm(expenses ~ ., data =  
insurance)
```

После того как модель построена, просто введите имя объекта модели, чтобы увидеть полученные бета-коэффициенты:

```
> ins_modelCall:lm(formula = expenses ~ ., data  
=  
insurance)Coefficients:      (Intercept)  
age      sexmale      -
```

11941.6	256.8	-	
131.4	bmi	children	smok
eryes	339.3	475.7	23
847.5	regionnorthwest	regionsoutheast	regionsout
hwest	-352.8	-1035.6	-
959.3			

Смысл этих коэффициентов регрессии довольно прост. Сдвиг — это прогнозируемое значение затрат, когда независимые переменные равны нулю. Однако во многих случаях сдвиг сам по себе имеет мало смысла, поскольку зачастую не бывает ситуации, при которой все признаки равны нулю. В нашем случае не существует людей с нулевым возрастом и нулевым коэффициентом ИМТ. Следовательно, у сдвига не существует интерпретации для реального мира. По этой причине на практике сдвиг часто игнорируется.

Бета-коэффициенты указывают на предполагаемое увеличение расходов при увеличении каждого признака на единицу при условии, что все остальные значения остаются постоянными. Например, ожидается, что с каждым годом медицинские расходы в среднем будут увеличиваться на 256,8 доллара при условии, что все остальные переменные не изменятся.

Аналогичным образом появление каждого ребенка приводит к дополнительным медицинским расходам, составляющим в среднем 475,7 доллара в год, а увеличение ИМТ на единицу вызывает увеличение ежегодных расходов на медицинское обслуживание в среднем на 339,3 доллара при прочих равных условиях.

Легко заметить, что, хотя в формуле модели указаны только шесть признаков, в результатах, кроме сдвига, указаны еще восемь коэффициентов. Это произошло потому, что функция `lm()` автоматически применяет фиктивное кодирование ко всем включенным в модель переменным, имеющим тип фактора.

Как было показано в главе 2, фиктивное кодирование позволяет рассматривать номинальный признак как числовой, создавая двоичную переменную для каждой категории этого признака. Фиктивная переменная принимает значение `1`, если наблюдение попадает в указанную категорию, и `0` — в противном случае. Например, признак `sex` имеет две категории: `male` и `female`. Он разбивается на две двоичные переменные, которым R присваивает имена `sexmale` и `sexfemale`. Для наблюдений, где `sex=male`, переменная `sexmale=1`, а `sexfemale=0`; и наоборот, если `sex=female`, то `sexmale=0` и `sexfemale=1`. Та же кодировка применяется к переменным с тремя и более категориями. Так, R разделяет признак `region` с четырьмя категориями на четыре фиктивные

переменные: `regionnorthwest`, `regionsoutheast`, `regionsouthwest` и `regionnortheast`.

При добавлении фиктивной переменной в регрессионную модель одна категория всегда является эталонной. Остальные оценки интерпретируются относительно эталонной. В нашей R-модели автоматически созданы переменные `sexfemale`, `smokerno` и `regionnortheast`, так что некурящие женщины, проживающие в северо-восточном регионе, становятся эталонной группой. Следовательно, у мужчин ежегодные медицинские расходы на 131,4 доллара меньше, чем у женщин, а курильщики тратят на медицину в среднем на 23 847,5 доллара в год больше, чем некурящие. Для всех оставшихся трех регионов в данной модели коэффициент является отрицательным. Это подразумевает, что эталонная группа — северо-восточный регион — имеет в среднем самые высокие расходы на медицину.



По умолчанию в R в качестве эталонного берется первый из уровней факторной переменной. Если вы предпочитаете сделать эталонным другой уровень, можно воспользоваться функцией `relevel()`, чтобы указать эталонную группу вручную. Для получения дополнительной информации используйте R-команду `relevel`.

В результатах, полученных с помощью линейной регрессионной модели, прослеживается логическая цепочка: преклонный возраст, курение и ожирение, как правило, вызывают проблемы со здоровьем, в то время как увеличение иждивенцев в семье может привести к более частым посещениям врача и обращению за профилактической помощью (прививки и ежегодные осмотры). Однако пока у нас нет представления о том, насколько хорошо модель соответствует данным. Ответ на этот вопрос мы получим в следующем разделе.

Шаг 4. Определение эффективности модели

Оценки параметров, которые мы получили, введя команду `ins_model`, показывают, как независимые переменные связаны с зависимой переменной, но не показывают, насколько хорошо модель соответствует данным. Чтобы оценить эффективность модели, мы можем использовать команду `summary()` для сохраненной модели:

```
> summary(ins_model)
```

В результате получим следующий вывод (цифры проставлены в иллюстративных целях):

```

Call:
lm(formula = expenses ~ ., data = insurance)

Residuals:
    Min       1Q   Median       3Q      Max    (1)
-11302.7 -2850.9  -979.6  1383.9 29981.7

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -11941.6     987.8  -12.089 < 2e-16 ***
age          256.8       11.9   21.586 < 2e-16 ***
sexmale     -131.3      332.9   -0.395 0.693255
bmi         339.3       28.6   11.864 < 2e-16 ***
children    475.7      137.8    3.452 0.000574 ***
smokeryes  23847.5     413.1   57.723 < 2e-16 ***
regionnorthwest -352.8    476.3   -0.741 0.458976
regionsoutheast -1035.6    478.7   -2.163 0.030685 *
regionsouthwest -959.3    477.9   -2.007 0.044921 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6062 on 1329 degrees of freedom
Multiple R-squared:  0.7509, Adjusted R-squared:  0.7494
F-statistic: 500.9 on 8 and 1329 DF, p-value: < 2.2e-16

```

На первый взгляд, результат, полученный с помощью функции `summary()`, может показаться невероятным, однако легко понять, на чем он основан. В таблице в пронумерованных метках видно, что есть три ключевых способа оценить эффективность нашей модели, то есть ее соответствие данным.

1. В разделе *Residuals* представлена сводная статистика ошибок наших прогнозов; некоторые из этих ошибок, по-видимому, весьма существенны. Поскольку остаточное значение равно разности истинного и прогнозируемого значений, максимальная ошибка составляет 29 981,7. Это предполагает, что как минимум для одного наблюдения модель занизила прогнозируемые расходы почти на 30 000 долларов. С другой стороны, 50 % ошибок относятся к значениям 1-го и 3-го квартала (первая и третья квартили), поэтому большинство прогнозов находятся в пределах от 2850,9 доллара больше истинного значения до 1383,9 доллара меньше истинного значения.

2. Для каждого вычисленного коэффициента регрессии определено *p*-значение, обозначаемое как `Pr(>|t|)`, которое представляет собой оценку вероятности того, что для данной оценки истинный коэффициент равен нулю. Небольшие *p*-значения предполагают, что истинный коэффициент вряд ли равен нулю, следовательно, маловероятно, что данный признак не связан с зависимой переменной. Обратите внимание, что некоторые из *p*-значений помечены звездочками (`***`), то есть имеют сноски. Эти сноски указывают на *уровень значимости* для данной оценки. Уровень значимости является порогом, выбранным до построения модели, который будет использоваться для обозначения «реальных», а не случайных результатов; *p*-значения ниже уровня значимости считаются *статистически значимыми*. Если

в модели мало таких значений, это может быть причиной для беспокойства, поскольку указывает на то, что используемые признаки не очень хорошо прогнозируют результат. В нашей модели есть несколько весьма значимых переменных, и они, похоже, логически связаны с результатом.

3. *Коэффициент детерминации* позволяет оценить, насколько хорошо модель в целом объясняет значения зависимой переменной. Эта величина похожа на коэффициент корреляции, так как чем ближе ее значение к 1,0, тем лучше модель объясняет данные. Поскольку в модели коэффициент детерминации равен 0,7494, это означает, что модель объясняет почти 75 % изменений зависимой переменной. Поскольку чем больше в модели признаков, тем лучше они объясняют изменения зависимой переменной, *коэффициент детерминации* представляет собой значение, скорректированное за счет наложения штрафов на модели с большим количеством независимых переменных. Это полезно для сравнения эффективности моделей с различным количеством объясняющих переменных.

На основании описанных трех показателей эффективности можно сказать, что наша модель работает весьма хорошо. Нередко регрессионные модели реальных данных имеют довольно низкий коэффициент детерминации; значение 0,75, по сути, довольно хороший результат. Немного беспокоит величина некоторых ошибок, но это не удивительно, учитывая характер данных о медицинских расходах. Однако, как вы увидите в следующем разделе, можно повысить эффективность модели, описав модель немного другим способом.

Шаг 5. Повышение эффективности модели

Как уже упоминалось, главное отличие регрессионного моделирования от других методов машинного обучения заключается в том, что регрессия обычно оставляет пользователю выбор признаков и спецификацию модели. Следовательно, если мы знаем о том, как признак связан с результатом, то можем использовать эту информацию в спецификации модели и, возможно, повысить ее эффективность.

Спецификация модели: добавление нелинейных зависимостей

В линейной регрессии предполагается, что связь между независимой и зависимой переменными является линейной, однако это не всегда так. Например, влияние возраста на медицинские расходы не может быть постоянным для всех возрастных значений; для пожилых групп населения лечение может стать слишком дорогим.

Напомню, что типичное уравнение регрессии имеет следующий вид:

$$y = \alpha + \beta_1 x$$

Чтобы учесть нелинейные зависимости, в уравнение регрессии можно добавить член более высокого порядка, представив модель в виде полинома. По сути, мы будем моделировать такие зависимости:

$$y = \alpha + \beta_1 x + \beta_2 x^2.$$

Разница между этими двумя моделями заключается в том, что вводится дополнительный бета-коэффициент, предназначенный для отражения влияния члена x^2 . Это позволяет представить влияние возраста как функцию квадрата возраста.

Для того чтобы добавить в модель нелинейную зависимость от возраста, нужно просто создать еще одну переменную:

```
> insurance$age2 <- insurance$age^2
```

Затем, когда мы построим улучшенную модель, добавим в формулу `lm()` как `age`, так и `age2`, в форме `expenses~age+age2`. Это позволит разделить в модели линейное и нелинейное влияние возраста на медицинские расходы.

Трансформация: преобразование числовой переменной в двоичный индикатор

Предположим, что влияние признака не является накопительным, скорее, оно начинает сказываться только после достижения определенного порогового значения. Например, ИМТ может никак не влиять на медицинские расходы для людей с нормальным весом, но сильно влиять на более высокие затраты для тучных людей (то есть для ИМТ со значениями 30 и более).

Можно смоделировать это воздействие, создав бинарную переменную индикатора ожирения, которая равна 1, если ИМТ превышает 30, и 0, если не превышает. Затем оценочный бета-коэффициент для этого бинарного признака будет указывать на среднее чистое влияние на медицинские расходы для тех лиц, чей ИМТ равен или превышает 30, относительно тех, у кого ИМТ менее 30.

Чтобы создать признак, можно воспользоваться функцией `ifelse()`, которая проверяет указанное условие для каждого элемента вектора и возвращает значение в зависимости от того, является ли это условие истинным или ложным. Для ИМТ, больше или равного 30, мы будем возвращать 1, в противном случае — 0:

```
> insurance$bmi30 <- ifelse(insurance$bmi >= 30, 1, 0)
```

Затем можно включить переменную `bmi30` в нашу улучшенную модель либо заменив ею исходную переменную `bmi`, либо в дополнение к ней, в зависимости от того, считаем ли мы, что воздействие ожирения дополняет отдельное линейное воздействие ИМТ. Без веских причин поступить иначе мы включим в нашу окончательную модель оба признака.



Если вы не уверены, стоит ли включать переменную в модель, обычной практикой является включить ее и проверить р-значение. Если переменная не является статистически значимой, то у вас появятся доказательства, дающие возможность впоследствии ее исключить.

Спецификация модели: добавление эффектов взаимодействия

До сих пор мы рассматривали только индивидуальный вклад каждого признака в результат. Но что, если отдельные признаки оказывают совместное воздействие на зависимую переменную? Например, курение и ожирение могут оказывать вредное воздействие по отдельности, но разумно предположить, что их совокупный эффект может оказаться хуже, чем от каждого из них в общей сложности.

Когда два признака имеют совокупный эффект, это называется *взаимодействием*. Если мы подозреваем, что две переменные взаимодействуют, можно проверить эту гипотезу, добавив их взаимодействие в модель. Эффекты взаимодействия задаются с использованием синтаксиса R-формулы. Чтобы задать взаимодействие между индикатором ожирения (`bmi30`) и индикатором курения (`smoker`), можно представить формулу в виде `expenses~bmi30*smoker`.

Оператор `*` — это сокращение, указывающее на то, что модель надо представить как `expenses~bmi30+smokeryes+bmi30:smokeryes`. Оператор двоеточия (`:`) в развернутом виде указывает на то, что `bmi30:smokeryes` представляет собой взаимодействие между двумя переменными. Обратите внимание, что расширенная форма также автоматически включает в себя и отдельные переменные `bmi30` и `smokeryes`, и их взаимодействие.



Взаимодействия никогда не должны включаться в модель без добавления каждой из взаимодействующих переменных. Если вы всегда создаете взаимодействия с помощью оператора `*`, это не проблема, поскольку R автоматически добавляет все необходимые компоненты.

Собираем все вместе: улучшенная регрессионная модель

Зная, каким образом медицинские расходы могут зависеть от характеристик пациента, мы разработали, как нам кажется, более точную формулу регрессии. В итоге мы внесли следующие улучшения:

- добавили нелинейную переменную для возраста;
- создали индикатор ожирения;
- ввели зависимость между ожирением и курением.

Мы будем обучать модель, как и раньше, используя функцию `lm()`, однако на этот раз добавим созданные переменные и эффект взаимодействия:

```
> ins_model2 <- lm(expenses ~ age + age2 +
children + bmi + sex
+ bmi30*smoker + region, data =
insurance)
```

Посмотрим, что получилось:

```
> summary(ins_model2)
```

Результат выглядит следующим образом:

```
Call:
lm(formula = expenses ~ age + age2 + children + bmi + sex + bmi30 *
smoker + region, data = insurance)
```

Residuals:

```
      Min       1Q   Median       3Q      Max
-17297.1 -1656.0 -1262.7  -727.8  24161.6
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	139.0053	1363.1359	0.102	0.918792
age	-32.6181	59.8250	-0.545	0.585690
age2	3.7307	0.7463	4.999	6.54e-07 ***
children	678.6017	105.8855	6.409	2.03e-10 ***
bmi	119.7715	34.2796	3.494	0.000492 ***
sexmale	-496.7690	244.3713	-2.033	0.042267 *
bmi30	-997.9355	422.9607	-2.359	0.018449 *
smokeryes	13404.5952	439.9591	30.468	< 2e-16 ***
regionnorthwest	-279.1661	349.2826	-0.799	0.424285
regionsoutheast	-828.0345	351.6484	-2.355	0.018682 *
regionsouthwest	-1222.1619	350.5314	-3.487	0.000505 ***
bmi30:smokeryes	19810.1534	604.6769	32.762	< 2e-16 ***

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 4445 on 1326 degrees of freedom
```

```
Multiple R-squared:  0.8664, Adjusted R-squared:  0.8653
```

```
F-statistic: 781.7 on 11 and 1326 DF, p-value: < 2.2e-16
```

Статистика соответствия модели помогает определить, повысилась ли эффективность регрессионной модели после внесенных изменений. По сравнению с первой моделью коэффициент детерминации повысился с 0,75 до 0,87.

Аналогично скорректированный коэффициент детерминации, который учитывает тот факт, что модель стала более сложной, повысился с 0,75 до 0,87. Теперь наша модель объясняет 87 % различий в стоимости лечения. Кроме того, наша теория о функциональной форме модели, похоже, подтверждается. Член более высокого порядка `age2` является статистически значимым, как и индикатор ожирения `bmi30`. Взаимосвязь между ожирением и курением имеет большое значение; кроме увеличения

затрат в случае курения на сумму более 13 404 долларов, курильщики, страдающие от ожирения, тратят еще 19 810 долларов в год. Это может свидетельствовать о том, что курение усугубляет заболевания, связанные с ожирением.



Строго говоря, регрессионное моделирование делает некоторые сильные предположения о данных. Для числового прогнозирования эти допущения не особенно важны, поскольку ценность модели не основана на том, действительно ли она отражает базовый процесс, — нас просто волнует точность прогнозов. Но если вы хотите сделать точные выводы на основании коэффициентов регрессионной модели, необходимо выполнить диагностические тесты, чтобы убедиться, что допущения регрессионной модели не были нарушены. Отличное введение в эту тему — книга Allison P.D. *Multiple Regression: A Primer*. Pine Forge Press, 1998.

Прогнозирование с помощью регрессионной модели

Изучив оценочные коэффициенты регрессии и статистику соответствия, мы можем использовать модель для прогнозирования расходов на медицинские страховки. Чтобы проиллюстрировать процесс прогнозирования, сначала применим модель к исходным тренировочным данным, используя функцию `predict()` следующим образом:

```
> insurance$pred <- predict(ins_model2,
insurance)
```

В результате прогнозы будут сохранены как новый вектор с именем `pred` во фрейме данных `insurance`. Затем можно вычислить корреляцию между прогнозируемой и фактической стоимостью страхования:

```
> cor(insurance$pred, insurance$expenses)[1]
0.9307999
```

Корреляция в 0,93 предполагает очень сильную линейную зависимость между прогнозируемыми и фактическими значениями. Это хороший знак — значит, модель очень точна! Полезно также представить результат в виде диаграммы рассеяния, которая приведена на рис. 6.9. Следующие R-команды строят отношения и добавляют единичную прямую — линию со сдвигом, равным 0, и наклоном, равным 1.

Параметры `col`, `lwd` и `lty` определяют цвет, толщину и тип линии соответственно:

```
> plot(insurance$pred, insurance$expenses)>
abline(a = 0, b = 1, col = "red", lwd = 3, lty =
2)
```

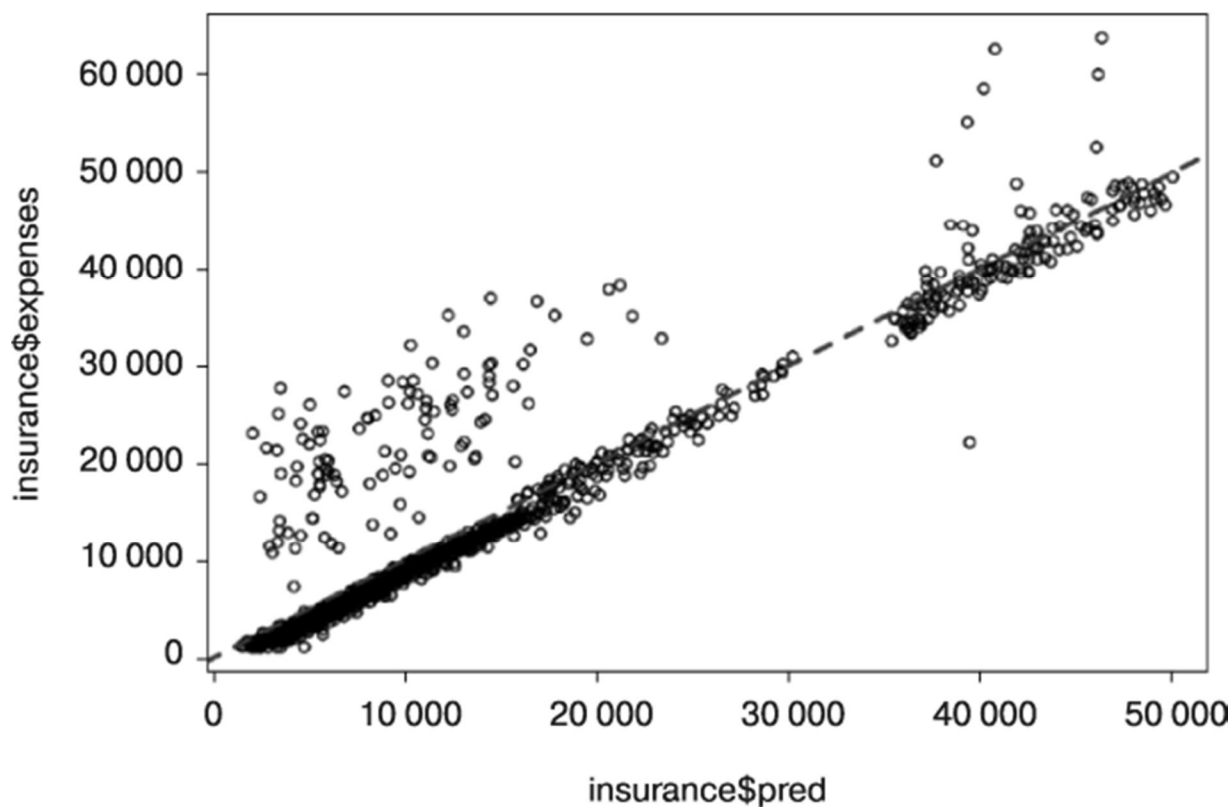


Рис. 6.9. На этой диаграмме рассеяния точки, попадающие на диагональную пунктирную линию $y = x$ или расположенные рядом с ней, указывают на прогнозы, очень близкие к фактическим значениям

Точки, расположенные выше диагональной прямой, соответствуют случаям, когда фактические расходы оказались больше, чем ожидалось, а точки, расположенные ниже этой линии, — случаям, когда расходы были меньше, чем ожидалось. Здесь мы видим, что небольшое количество пациентов с гораздо большими, чем ожидалось, медицинскими расходами уравнивается большим количеством пациентов с расходами, немного меньше ожидаемых.

Теперь предположим, что мы хотим спрогнозировать расходы для новых участников программы страхования. Для этого нужно предоставить функции `predict()` фрейм данных с информацией о предполагаемых пациентах. Если пациентов много, можно создать файл в формате электронной таблицы CSV для загрузки в R; если же пациентов немного, то можно просто создать фрейм данных в самой функции `predict()`. Например, для того, чтобы оценить расходы на страхование 30-летнего мужчины, некурящего, но имеющего избыточный вес, проживающего с двумя детьми на северо-востоке США, можно ввести следующее:

```
> predict(ins_model2, data.frame(age =
30, age2 = 30^2, children =
2, bmi = 30, sex = "male",
bmi30 = 1, smoker = "no",
region = "northeast")) 15973.774
```

Используя это значение, страховая компания, возможно, решит установить цену для данной демографической группы на уровне около 6000 долларов в год, или 500 долларов в месяц, чтобы гарантировать безубыточность. Чтобы сравнить этот показатель для женщины с такими

же характеристиками, можно использовать функцию `predict()` почти таким же образом:

```
> predict(ins_model2, data.frame(age = 30, age2 = 30^2, children = 2, bmi = 30, sex = "female", bmi30 = 1, smoker = "no", region = "northeast")) 16470.543
```

Обратите внимание, что разность между этими двумя значениями составляет $5973,774 - 6470,543 = -496,769$ и равна расчетному коэффициенту регрессионной модели для `sexmale`. По оценкам, в среднем мужчины при прочих равных условиях тратят на медицину примерно на 496 долларов в год меньше, чем женщины.

Это свидетельствует о том, что прогнозируемые расходы являются суммой каждого из коэффициентов регрессии, умноженного на соответствующее значение во фрейме данных прогнозирования. Например, используя коэффициент регрессии модели 678,6017 для количества детей, можно предположить, что те, у кого нет детей, тратят на медицинское обслуживание на $2 \times 678,6017 = 1357,203$ доллара меньше:

```
> predict(ins_model2, data.frame(age = 30, age2 = 30^2, children = 0, bmi = 30, sex = "female", bmi30 = 1, smoker = "no", region = "northeast")) 15113.34> 6470.543 - 5113.34[1] 1357.203
```

Выполнив аналогичные операции еще для нескольких клиентских сегментов, страховая компания может разработать выгодную структуру ценообразования для различных демографических показателей.



Экспорт коэффициентов регрессионной модели позволяет построить собственную функцию прогнозирования. Одним из возможных вариантов использования такой функции было бы внедрение регрессионной модели в базу данных клиентов для прогнозирования в реальном времени.

Регрессионные деревья и деревья моделей

Вспомним, что говорилось в главе 5: дерево решений строит модель, очень похожую на потоковую диаграмму, где узлы решений, концевые узлы и ветви определяют последовательность решений, которые принимаются для классификации примеров. Такие деревья можно также использовать для числового прогнозирования, внося лишь небольшие изменения в алгоритм построения деревьев. В этом разделе рассмотрим, чем деревья для

числового прогнозирования отличаются от деревьев, используемых для классификации.

Деревья для числового прогнозирования делятся на две категории. Первая — *регрессионные деревья* — была введена в 1980-х годах как часть фундаментального алгоритма построения *дерева классификации и регрессии* (Classification and Regression Tree, CART). Несмотря на название, в регрессионных деревьях, как уже говорилось в этой главе, не используются методы линейной регрессии; вместо этого они делают прогнозы на основе среднего значения для имеющихся примеров, которые в итоге позволяют достичь конечного узла.



Алгоритм CART подробно описан в книге: Breiman L., Friedman J.H., Stone C.J., Olshen R.A. Classification and Regression Trees. Chapman and Hall, 1984. Вторая категория деревьев для числового прогнозирования — *деревья моделей*. Они появились на несколько лет позже, чем регрессионные деревья, и они менее известны, однако, возможно, более эффективны. Деревья моделей строятся примерно так же, как и регрессионные деревья, однако для каждого конечного узла строится модель множественной линейной регрессии на основе примеров, дошедших до этого узла. В зависимости от количества конечных узлов, дерево моделей может строить десятки или даже сотни таких моделей. Из-за этого деревья моделей труднее исследовать, чем аналогичные регрессионные деревья, но их преимущество заключается в том, что деревья моделей могут привести к построению более точной модели.



Самый первый алгоритм дерева моделей, M5, описан в публикации: Quinlan J.R. Learning with continuous classes // Proceedings of the 5th Australian Joint Conference on Artificial Intelligence, 1992. P. 343–348.

Дополнение деревьев регрессией

Деревья, способные выполнять числовое прогнозирование, представляют собой привлекательную альтернативу регрессионному моделированию, которую, однако, часто упускают из виду. В табл. 6.2 представлены преимущества и недостатки регрессионных деревьев и деревьев моделей по сравнению с более распространенными методами регрессии.

Таблица 6.2

Преимущества	Недостатки
<p>Сочетают в себе преимущества деревьев решений с возможностью моделирования числовых данных.</p> <p>Не требуют от пользователя заранее выбрать модель.</p> <p>Признаки выбираются автоматически, что позволяет использовать очень большое количество признаков.</p> <p>Для некоторых типов данных подходят намного лучше, чем линейная регрессия.</p> <p>Для интерпретации модели не требуется знание статистики</p>	<p>Не так известны, как линейная регрессия.</p> <p>Требуют большого количества тренировочных данных.</p> <p>Сложно определить общее результирующее влияние отдельных признаков на результат.</p> <p>Большие деревья могут оказаться более трудными для интерпретации, чем регрессионная модель</p>

При решении задач числового прогнозирования обычно выбирают традиционные регрессионные методы, однако иногда числовые деревья решений имеют явные преимущества. Например, деревья решений могут лучше подойти для задач с большим количеством признаков или множеством сложных нелинейных зависимостей между признаками и результатом; такие ситуации создают проблемы для регрессии. Регрессионное моделирование также делает предположения о данных, которые на практике часто нарушаются; в случае с деревьями этого не происходит.

Деревья числового прогнозирования строятся во многом так же, как и для классификации. Начиная с корневого узла, данные разделяются по стратегии «разделяй и властвуй» в соответствии с признаком, который после деления приведет к наибольшей однородности результатов. В классификационных деревьях, как мы помним, однородность измеряется энтропией. Для числовых данных это неприменимо. Вместо этого для числовых деревьев решений однородность измеряется статистическими параметрами, такими как дисперсия, стандартное отклонение или абсолютное отклонение от среднего.

Одним из наиболее распространенных критериев деления является *уменьшение стандартного отклонения* (Standard Deviation Reduction, SDR), которое вычисляется по следующей формуле:

$$SDR = sd(T) - \sum_i \frac{|T_i|}{|T|} \times sd(T_i)$$

В этой формуле функция $sd(T)$ — стандартное отклонение значений множества T , а T_1, T_2, \dots, T_n — множества значений, полученных в результате деления по заданному признаку. $|T|$ означает число наблюдений в множестве T . В сущности, эта формула позволяет вычислить уменьшение стандартного отклонения путем сравнения стандартного отклонения до деления со взвешенным стандартным отклонением после деления.

В качестве примера рассмотрим следующий случай: дерево принимает решение, следует ли выполнить деление по двоичному признаку А или по двоичному признаку Б (рис. 6.10).

Исходные данные	1	1	1	2	2	3	4	5	5	6	6	7	7	7	7
Разделение по признаку А	1	1	1	2	2	3	4	5	5	6	6	7	7	7	7
Разделение по признаку Б	1	1	1	2	2	3	4	5	5	6	6	7	7	7	7
	T_1							T_2							

Рис. 6.10. Алгоритм рассматривает варианты разделения множества по признакам А и Б, в результате которых создаются группы T_1 и T_2

На основе групп, которые получились бы в результате предложенных разделений, можно вычислить SDR для А и Б следующим образом. Используемая здесь функция `length()` возвращает количество элементов в векторе. Обратите внимание, что общая группа называется `tee`, чтобы избежать перезаписи встроенных R-функций `T()` и `t()`.

```
> tee <- c(1, 1, 1, 2, 2, 3, 4, 5, 5, 6, 6, 7,
7, 7, 7)> at1 <- c(1, 1, 1, 2, 2, 3, 4, 5, 5)>
at2 <- c(6, 6, 7, 7, 7, 7)> bt1 <- c(1, 1, 1, 2,
2, 3, 4)> bt2 <- c(5, 5, 6, 6, 7, 7, 7, 7)> sdr_a
<- sd(tee) - (length(at1) / length(tee) * sd(at1)
+
length(at2) / length(tee) *
sd(at2))> sdr_b <- sd(tee) - (length(bt1) /
length(tee) * sd(bt1) +
length(bt2) /
length(tee) * sd(bt2))
```

Сравним SDR для А и Б:

```
> sdr_a[1] 1.202815> sdr_b[1] 1.392751
```

SDR для разделения по признаку А составляет примерно 1,2, а для разделения по признаку Б — 1,4. Поскольку уменьшение стандартного отклонения для Б больше, чем для А, то дерево решений сначала будет использовать признак Б, так как в результате получаются более однородные подмножества, чем при разделении по признаку А.

Предположим, что в этот момент дерево перестало расти и ограничилось этим единственным разделением. На этом работа регрессионного дерева заканчивается. Можно сделать прогнозы для новых примеров в зависимости от того, попадает ли значение признака Б для данного примера в группу T_1 или T_2 . Если пример попадает в T_1 , то модель будет давать прогноз `mean(bt1)=2`, а если в T_2 — то `mean(bt2)=6.25`.

Дерево моделей, напротив, пошло бы дальше. Используя семь обучающих примеров, попадающих в группу T_1 , и восемь примеров, попадающих в группу T_2 , дерево моделей может построить линейную регрессионную модель результата по признаку А. Обратите внимание: признак Б не участвует в построении регрессионной модели, поскольку все примеры в списке имеют одинаковое значение Б — по этому признаку они были помещены в группу T_1 или T_2 . После этого дерево моделей может делать прогнозы для новых примеров, используя одну из двух линейных моделей.

Для того чтобы проиллюстрировать различия между этими двумя подходами, рассмотрим пример из реальной практики.

Пример: оценка качества вина с помощью регрессионного дерева и дерева моделей

Виноделие — сложный и конкурентный бизнес с высоким потенциалом для получения прибыли. Рентабельность винодельни зависит от множества факторов. Как и для любого сельскохозяйственного продукта, на качество вина влияют такие параметры, как погодные условия и среда произрастания. На вкусовые качества вина — в лучшую или худшую сторону — также влияют процесс производства и розлив. На восприятие продукта потребителем оказывает воздействие даже способ сбыта — от дизайна бутылки до цены.

Как следствие, винодельческая промышленность вкладывает значительные средства в методы сбора данных и машинного обучения, которые могут помочь в принятии решений в области виноделия. Например, машинное обучение использовалось для выявления ключевых различий в химическом составе вин из разных регионов, а также для выявления химических веществ, которые делают вино более сладким.

Совсем недавно машинное обучение использовалось для того, чтобы помочь в такой трудной задаче, как оценка качества вина. Обзор, составленный известным сомелье, часто помогает определить, какое место на полке магазина займет товар, хотя даже самые лучшие эксперты бывают непоследовательными при оценке вина вслепую.

В этом примере мы будем использовать регрессионные деревья и деревья моделей, чтобы построить систему, способную имитировать экспертные оценки вина. Поскольку деревья приводят к понятной модели, это может помочь виноделам выделить ключевые факторы, способствующие производству вин с более высоким рейтингом. Что еще важнее, наша система не будет подвержена человеческому фактору в процессе дегустации, такому как настроение или нарушения в работе рецепторов. Таким образом, компьютерное тестирование вина может привести к получению лучшего продукта, а также к более объективным, последовательным и справедливым оценкам.

Шаг 1. Сбор данных

Для разработки модели рейтинга вин мы будем использовать данные, безвозмездно переданные в репозиторий машинного обучения *UCI Machine Learning Repository* (<http://archive.ics.uci.edu/ml>) П. Кортесом (P. Cortez), А. Сердейрой (A. Cerdeira), Ф. Алмейдой (F. Almeida), Т. Матосом (T. Matos) и Дж. Рейсом (J. Reis). Этот набор данных включает в себя образцы красных и белых вин «Винью-Верде» из Португалии — одной из ведущих

стран — производителей вина в мире. Поскольку факторы, влияющие на высокую оценку вина, разные для красных и белых сортов, в этом анализе мы рассмотрим только более популярные белые вина.



Для того чтобы выполнить этот пример, загрузите вместе с материалами к книге файл `whitewines.csv` и сохраните его в рабочем каталоге R. Там также есть файл `redwines.csv` — на тот случай, если вы захотите исследовать эти данные самостоятельно.

Данные о белых винах включают в себя информацию об 11 химических свойствах 4898 образцов вин. Для каждого вина были проведены лабораторные исследования для определения кислотности, содержания сахара, хлоридов, серы, спирта, а также уровня pH и плотности. Затем образцы оценивались методом слепой дегустации группами из не менее трех экспертов по шкале качества от нуля (очень плохо) до десяти (отлично). В случае если эксперты не приходили к единому мнению, использовалось медианное значение.

Исследование, проведенное Кортесом, оценивало возможности трех методов машинного обучения к моделированию данных о винах: множественной регрессии, искусственных нейронных сетей и метода опорных векторов. Множественная регрессия уже рассматривалась в этой главе; о нейронных сетях и методах опорных векторов вы узнаете из главы 7. Исследование показало, что метод опорных векторов дает значительно лучшие результаты, чем линейная регрессионная модель. Однако, в отличие от регрессионной модели, модель метода опорных векторов трудно интерпретировать. Используя регрессионные деревья и деревья моделей, можно улучшить результаты регрессии, в то же время получая модель, которую легко понять.



Подробнее об описанном здесь исследовании вин читайте в статье: Cortez P., Cerdeira A., Almeida F., Matos T., Reis J. Modeling wine preferences by data mining from physicochemical properties. *Decision Support Systems*, 2009. Vol. 47. P. 547–553.

Шаг 2. Исследование и подготовка данных

Как обычно, для загрузки данных в среду R мы воспользуемся функцией `read.csv()`. Поскольку все признаки являются числовыми, мы можем смело игнорировать параметр `stringsAsFactors`:

```
> wine <- read.csv("whitewines.csv")
```

Данные о винах включают в себя 11 признаков и вывод о качестве, а именно:

```

> str(wine)'data.frame':   4898 obs. of  12
variables:$ fixed.acidity      : num  6.7 5.7
5.9 5.3 6.4 7 7.9 ...$ volatile.acidity      :
num  0.62 0.22 0.19 0.47 0.29 0.12 ...$
citric.acid      : num  0.24 0.2 0.26 0.1 0.21
0.41 ...$ residual.sugar      : num  1.1 16 7.4
1.3 9.65 0.9 ...$ chlorides      :
num  0.039 0.044 0.034 0.036 0.041 ...$
free.sulfur.dioxide : num  6 41 33 11 36 22 33 17
34 40 ...$ total.sulfur.dioxide: num  62 113 123
74 119 95 152 ...$ density      :
num  0.993 0.999 0.995 0.991 0.993 ...$
pH      : num  3.41 3.22 3.49 3.48
2.99 3.25 ...$ sulphates      : num  0.32
0.46 0.42 0.54 0.34 0.43 ...$
alcohol      : num  10.4 8.9 10.1 11.2
10.9 ...$ quality      : int  5 6 6 4 6 6
6 6 6 7 ...

```

По сравнению с другими типами моделей машинного обучения, одним из преимуществ деревьев является то, что они способны обрабатывать различные типы данных без предварительной обработки. Это означает, что не нужно нормализовывать или стандартизировать признаки.

Однако необходимо приложить усилия для исследования распределения результирующей переменной, чтобы затем можно было оценить эффективность модели. Предположим, что различия в качестве вин были очень незначительными или что распределение соответствует бимодальному распределению: либо очень хорошее, либо очень плохое. Это может повлиять на построение модели. Для того чтобы проверить существование таких крайностей, можно рассмотреть распределение качества вин, построив гистограмму:

```

> hist(wine$quality)

```

В результате получим следующий график (рис. 6.11).

Как видим, показатели качества вин вполне соответствуют нормальному распределению, образуя гистограмму в виде колокола, с центральным значением, примерно равным шести. Интуитивно это имеет смысл, поскольку большинство вин среднего качества; лишь немногие из них исключительно хороши или плохи. Здесь результаты не показаны, однако также полезно проверить результат выполнения функции `summary(wine)` на предмет выбросов и других потенциальных проблем данных. Несмотря на то что деревья весьма устойчивы к ненадежным данным, всегда целесообразно проверять данные

на наличие серьезных проблем. Сейчас мы предположим, что данные надежны.

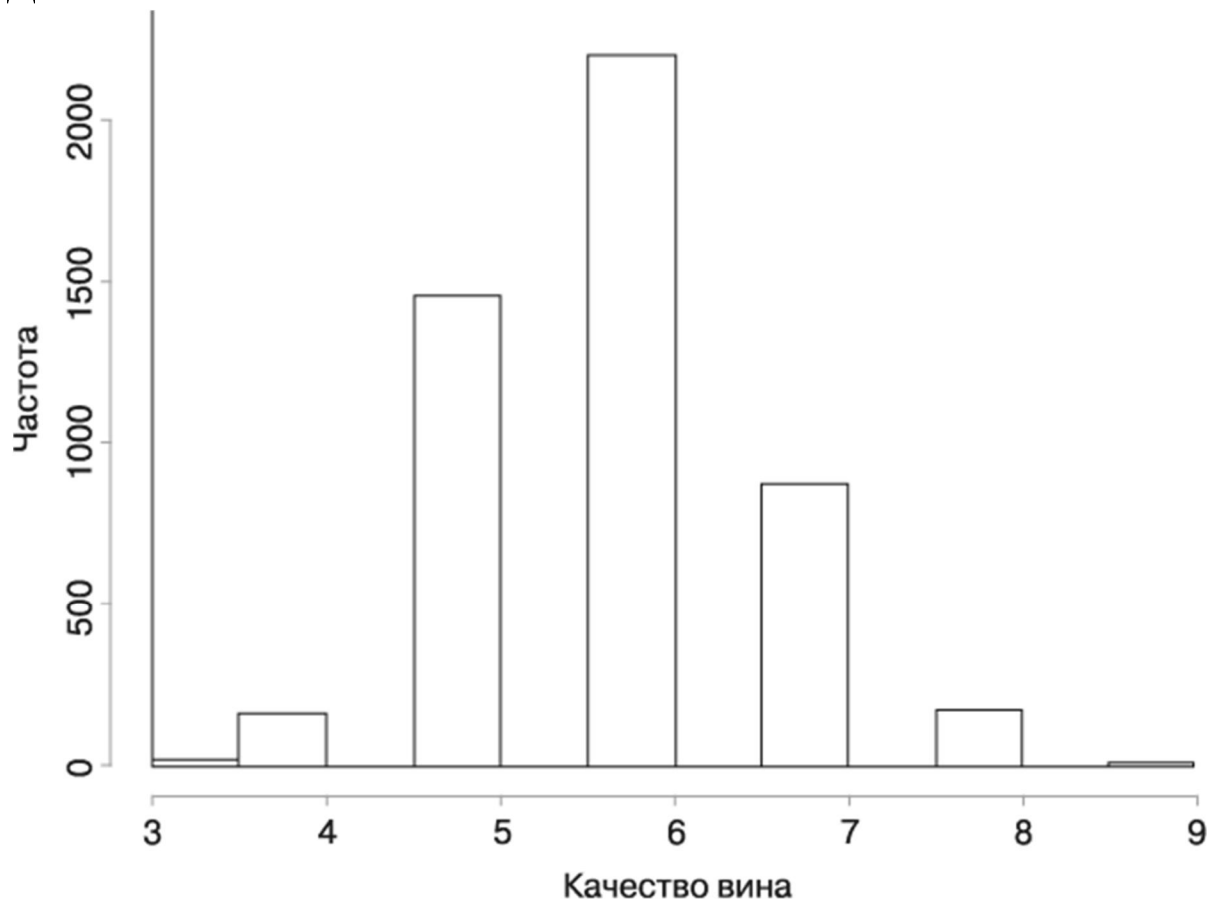


Рис. 6.11. Распределение рейтинга качества белых вин

Теперь остался последний шаг: разделить набор данных на наборы для обучения и тестирования. Поскольку набор данных о винах упорядочен случайным образом, мы можем разбить его на два набора смежных строк:

```
> wine_train <- wine[1:3750, ] > wine_test <-  
wine[3751:4898, ]
```

Для того чтобы отразить условия, использованные Кортесом, мы разделили набор данных так, что тренировочные данные составили 75 %, а тестовые — 25 % исходного набора. С помощью тестовых данных оценим эффективность наших моделей деревьев, чтобы увидеть, удалось ли нам получить результаты, сопоставимые с результатами предыдущих исследований.

Шаг 3. Обучение модели на данных

Начнем с обучения модели регрессионного дерева. Для моделирования регрессионных деревьев можно использовать практически любую реализацию деревьев решений, однако наиболее точная реализация регрессионных деревьев в том виде, как они описаны командой CART, предлагается в пакете `rpart` (от *recursive partitioning* — рекурсивное разбиение). Как и классическая R-реализация CART, пакет `rpart` хорошо документирован и поддерживается функциями визуализации и оценки моделей `rpart`.

Установите пакет `rpart` с помощью команды `install.packages("rpart")`. Затем его можно загрузить в сеанс R с помощью команды `library(rpart)`. Следующий

синтаксис отвечает обучению дерева с использованием настроек по умолчанию, которые обычно довольно хорошо работают. Если понадобится более тонкая настройка, обратитесь к документации по параметрам управления, введя команду `?rpart.control`.

Синтаксис регрессионной модели

Построение модели:

```
m <- rpart(dv ~ iv, data = mydata)
```

`dv` – зависимая переменная из фрейма данных `mydata`, которую нужно смоделировать;

`iv` – R-формула, определяющая независимые переменные во фрейме данных `mydata`, которые следует использовать в модели;

`data` – фрейм данных, которому принадлежат переменные `dv` и `iv`.

Функция возвращает объект модели регрессионного дерева, который может быть использован для прогнозирования.

Прогнозирование:

```
p <- predict(m, test, type = "vector")
```

`m` – модель, обученная с помощью функции `rpart()`;

`test` – фрейм данных, содержащий тестовые данные с теми же признаками, что и тренировочные данные, использованные при построении модели;

`type` – тип возвращаемого прогноза: `"vector"` (для прогнозирования числовых значений), `"class"` (для прогнозирования классов) или `"vector"` (для прогнозирования вероятности принадлежности к классу).

Функция возвращает вектор прогнозируемых значений в зависимости от значения параметра `type`.

Пример:

```
wine_model <- rpart(quality ~ alcohol + sulfates,
                    data = wine_train)
wine_predictions <- predict(wine_model, wine_test)
```

Можно указать `quality` в качестве выходной переменной и использовать точечную нотацию, чтобы задействовать остальные столбцы во фрейме данных `wine_train` в качестве предикторов. Полученный в результате объект модели регрессионного дерева называется `m.rpart`, чтобы отличать его от дерева моделей, которое мы рассмотрим позже:

```
> m.rpart <- rpart(quality ~ ., data = wine_train)
```

Для получения базовой информации о дереве просто введите имя объекта модели:

```
> m.rpartn= 3750node), split, n, deviance,
yval      * denotes terminal node) root 3750
2945.53200 5.870933      2) alcohol < 10.85 2372
1418.86100 5.604975      4)
volatile.acidity >= 0.2275 1611 821.30730
5.432030      8) volatile.acidity >= 0.3025 688
```

```

278.97670 5.255814 *          9) volatile.acidity<
0.3025 923 505.04230 5.563380 *          5)
volatile.acidity< 0.2275 761 447.36400 5.971091
*          3) alcohol>=10.85 1378 1070.08200
6.328737          6) free.sulfur.dioxide< 10.5 84
95.55952 5.369048 *          7)
free.sulfur.dioxide>=10.5 1294 892.13600
6.391036          14) alcohol< 11.76667 629 430.11130
6.173291          28) volatile.acidity>=0.465 11
10.72727 4.545455 *          29) volatile.acidity<
0.465 618 389.71680 6.202265 *          15)
alcohol>=11.76667 665 403.99400 6.596992 *

```

Для каждого узла в дереве указано количество примеров, достигающих точки принятия решения. Например, все 3750 примеров начинаются с корневого узла, из которых для 2372 `alcohol<10.85`, а для 1378 — `alcohol>=10.85`. Поскольку признак `alcohol` стал первым используемым в дереве, это самый важный показатель качества вина.

Узлы, помеченные знаком `*`, являются терминальными, или концевыми, узлами. Это означает, что они приводят к прогнозу (представленному здесь как `yval`). Например, для узла 5 `yval` = 5,971091. Поэтому, когда дерево используется для прогнозов, все образцы вин, для которых `alcohol<10.85` и `volatile.acidity<0,2275`, получают прогнозируемое значение `quality`, равное 5,97.

Более подробное описание соответствия дерева, включая среднеквадратичную ошибку для каждого из узлов и общую меру важности признаков, можно получить с помощью команды `summary(m.rpart)`.

Визуализация деревьев решений

Чтобы понять, как работает дерево, достаточно предыдущего вывода `m.rpart`, однако зачастую проще использовать визуализацию. В состав пакета `rpart.plot`, разработанного Стивеном Милборроу (Stephen Milborrow), входит простая функция, которая строит деревья решений, удобные для графического представления.



Для получения дополнительной информации о функции `rpart.plot`, включая примеры типов диаграмм для деревьев решений, которые может создавать эта функция, обратитесь на мой сайт по адресу <http://www.milbo.org/rpart-plot/>.

После установки пакета с помощью команды `install.packages("rpart.plot")` функция `rpart.plot()` строит древовидную диаграмму на основании любого объекта модели `rpart`. Следующие команды строят такое же регрессионное дерево, как и то, которое мы создали ранее:

```
> library(rpart.plot) > rpart.plot(m.rpart, digits = 3)
```

В результате получим древовидную диаграмму, приведенную на рис. 6.12.

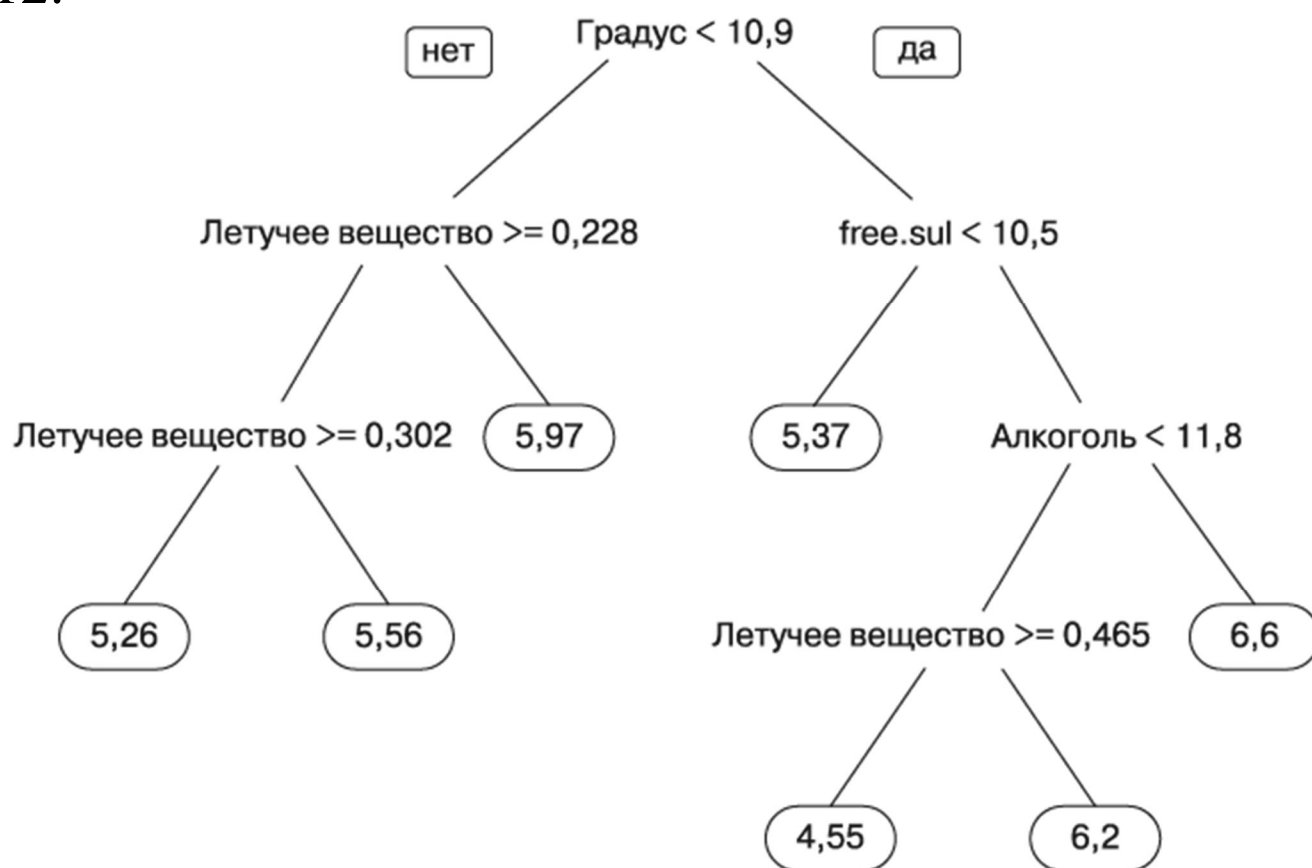


Рис. 6.12. Визуализация модели регрессионного дерева для определения качества вин

Кроме параметра `digits`, который определяет количество значащих цифр, используемых на диаграмме, можно настроить многие другие аспекты визуализации. В следующей команде использованы лишь некоторые полезные параметры:

```
> rpart.plot(m.rpart, digits = 4, fallen.leaves = TRUE, type = 3, extra = 101)
```

Параметр `fallen.leaves` выравнивает концевые узлы внизу графика, а параметры `type` и `extra` определяют способ маркировки решений и узлов. Числа `3` и `101` соответствуют конкретным форматам стилей, полный список которых вы найдете в документации к команде или экспериментируя с различными номерами.

Результатом этих изменений является древовидная диаграмма, представленная на рис. 6.13, которая выглядит совершенно иначе.

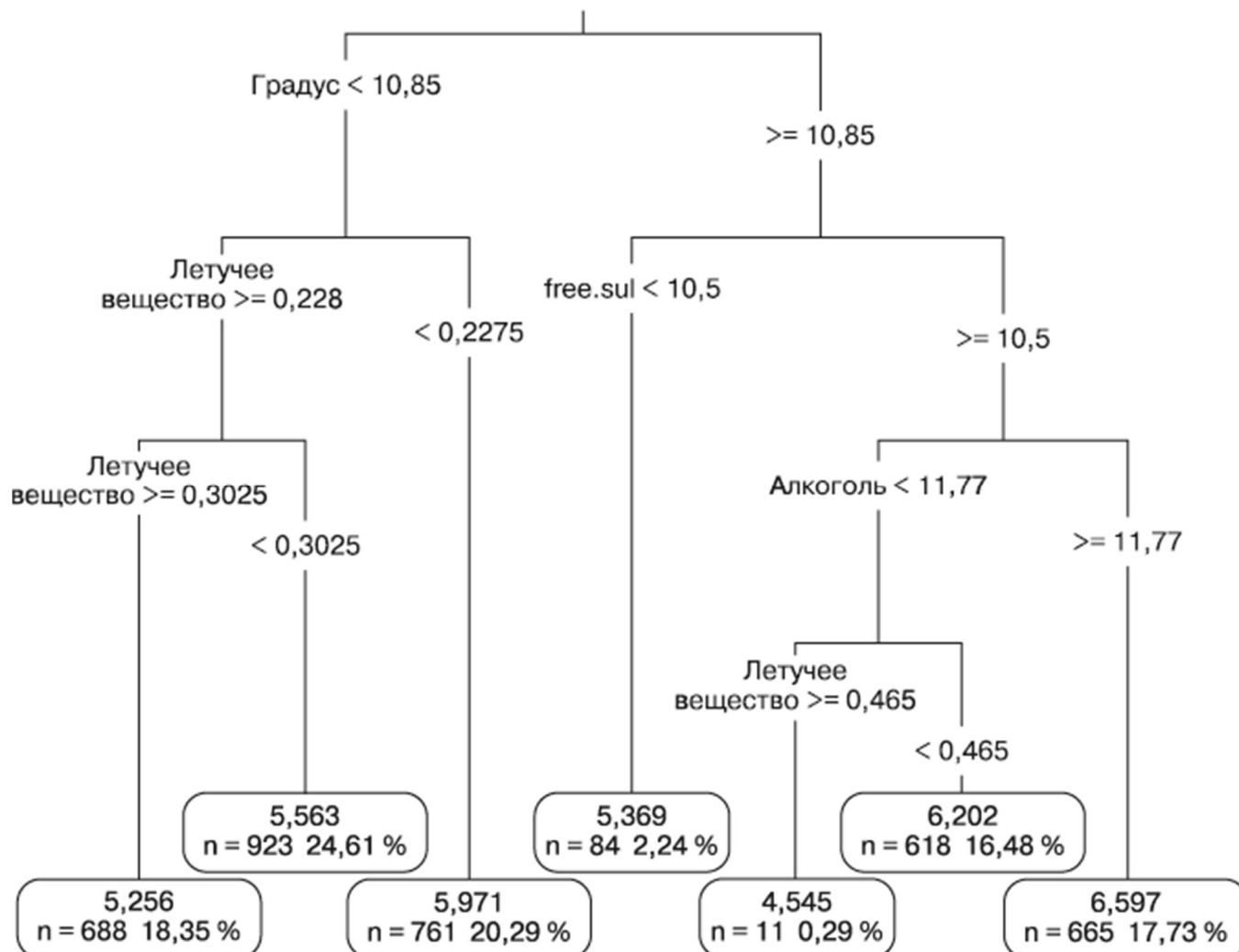


Рис. 6.13. Изменение параметров функции построения диаграммы позволяет настроить визуализацию дерева

Подобные визуализации помогают распространять результаты построения регрессионных деревьев, поскольку их легко понять, даже не имея математического образования. В обоих случаях числа, показанные в конечных узлах, являются спрогнозированными значениями для примеров, достигающих этого узла. Таким образом, если показать диаграмму производителям вина, это может помочь им определить ключевые факторы, которые позволяют производить вина с более высоким рейтингом.

Шаг 4. Определение эффективности модели

Для того чтобы использовать модель регрессионного дерева для прогнозирования на основе тестовых данных, мы воспользуемся функцией `predict()`. По умолчанию она возвращает приблизительное числовое значение результирующей переменной, которое мы сохраним в векторе `p.rpart`:

```
> p.rpart <- predict(m.rpart, wine_test)
```

Достаточно взглянуть на сводную статистику полученных прогнозов, чтобы заметить потенциальную проблему: у прогнозов гораздо более узкий диапазон, чем у реальных значений:

```
> summary(p.rpart)  Min.  1st
Qu.  Median  Mean  3rd
Qu.    Max. 4.545  5.563  5.971  5.893  6.202
6.597> summary(wine_test$quality)  Min.  1st
```

Qu.	Median	Mean	3rd			
Qu.	Max.	3.000	5.000	6.000	5.901	6.000
		9.000				

Этот результат свидетельствует о том, что модель неправильно определяет экстремальные случаи, в частности лучшие и худшие вина. Однако между первым и третьим квартилем она может работать хорошо.

Корреляция между прогнозируемыми и реальными значениями качества обеспечивает простой способ измерения эффективности модели. Напомню, что для измерения отношения между двумя векторами равной длины можно использовать функцию `cor()`. Мы воспользуемся ею, чтобы сравнить, насколько хорошо прогнозы соответствуют реальным значениям:

```
> cor(p.rpart, wine_test$quality)[1] 0.5369525
```

Корреляция 0,54, безусловно, приемлема. Однако она лишь показывает, насколько сильно прогнозы связаны с реальным значением, но не то, насколько далеки прогнозы от истинных значений.

Измерение эффективности с помощью средней абсолютной ошибки
Еще один способ оценить эффективность модели — определить, насколько далек в среднем ее прогноз от истинного значения. Такое измерение называется *средней абсолютной ошибкой* (Mean Absolute Error, MAE). Уравнение для вычисления MAE выглядит следующим образом, где n означает количество прогнозов, a^e_i — ошибку для i -го прогноза:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |e_i|$$

Как следует из названия, это уравнение принимает среднее значение абсолютной величины ошибок. Поскольку ошибка — всего лишь разность между прогнозируемым и фактическим значением, можно построить простую функцию `MAE()` следующим образом:

```
> MAE <- function(actual, predicted)
{ mean(abs(actual - predicted)) }
```

MAE для наших прогнозов выглядит так:

```
> MAE(p.rpart, wine_test$quality)[1] 0.5872652
```

Это означает, что в среднем разница между прогнозами нашей модели и истинным показателем качества составила примерно 0,59. При шкале качества от нуля до десяти это говорит о том, что данная модель работает достаточно хорошо.

Однако напомню, что большинство вин не были ни очень хорошими, ни очень плохими; типичный показатель качества колебался около 5–6. Следовательно, классификатор, который бы не делал ничего, кроме вычисления среднего значения, в соответствии с этой метрикой тоже работал бы достаточно хорошо.

Средняя оценка качества для тренировочных данных выглядит следующим образом:

```
> mean(wine_train$quality)[1] 5.870933
```

Если бы мы давали прогноз 5,87 для всех образцов вин, то получили бы среднюю абсолютную ошибку всего лишь около 0,67:

```
> MAE(5.87, wine_test$quality)[1] 0.6722474
```

Наше регрессионное дерево (MAE = 0,59) в среднем ближе к истинному показателю качества, чем условное среднее (MAE = 0,67), но не намного. Для сравнения, Кортес сообщил о том, что для модели нейронной сети значение MAE составило 0,58, а для метода опорных векторов — 0,45. Это говорит о том, что можно еще улучшить алгоритм.

Шаг 5. Повышение эффективности модели

Для того чтобы повысить эффективность нашего алгоритма обучения, мы применим алгоритм дерева моделей, который является более сложной реализацией деревьев для числового прогнозирования. Напомню, что деревья моделей расширяют регрессионные деревья, заменяя концевые узлы регрессионными моделями. Это часто дает более точные результаты, чем регрессионные деревья, которые для прогнозирования в концевых узлах используют только одно числовое значение.

Самое современное из деревьев моделей — это алгоритм *Cubist*, который, в свою очередь, является усовершенствованным алгоритмом дерева моделей M5. Оба алгоритма были описаны Дж. Р. Куинланом (J. R. Quinlan) в начале 1990-х годов. Детали реализации этих алгоритмов выходят за рамки данной книги, однако отмечу, что алгоритм *Cubist* включает в себя построение дерева решений, создание правил принятия решений на основе ветвей дерева и построение регрессионной модели в каждом из концевых узлов. Для улучшения качества прогнозов и сглаживания во всем диапазоне прогнозируемых значений используется дополнительная эвристика, такая как сокращение и бустинг.



Подробнее об алгоритмах *Cubist* и M5 читайте в статье: Quinlan J.R. Learning With Continuous Classes // Proceedings of the 5th Australian Joint Conference on Artificial Intelligence, 1992. P. 343–348, а также в статье: Quinlan J.R. Combining Instance-Based and Model-Based Learning // Proceedings of the Tenth International Conference on Machine Learning, 1993. P. 236–243.

Алгоритм *Cubist* в R реализован в составе пакета `Cubist` в виде входящей в него функции `cubist()`. Синтаксис этой функции представлен в следующей таблице.

Синтаксис деревьев моделей

Построение модели:

```
m <- cubist(train, class)
```

`train` – фрейм данных или матрица с тренировочными данными;

`class` – факторный вектор, каждая строка которого содержит класс для тренировочных данных.

Функция возвращает объект дерева `cubist`-модели, которую можно использовать для прогнозирования.

Прогнозирование:

```
p <- predict(m, test)
```

`m` – модель, обученная с помощью функции `cubist()`;

`test` – фрейм данных, содержащий тестовые данные с теми же признаками, что и тренировочные данные, использованные при построении модели.

Функция возвращает вектор числовых прогнозируемых значений.

Пример:

```
wine_model <- cubist(wine_train, wine_quality)
wine_predictions <- predict(wine_model, wine_test)
```

При построении дерева модели `Cubist` мы будем использовать синтаксис, немного отличающийся от того, который использовался для регрессионного дерева, так как функция `cubist()` не воспринимает R-формулы. Вместо этого нам нужно указать столбцы фрейма данных, соответствующие независимым переменным `x` и зависимой переменной `y`. Поскольку прогнозируемое качество вина находится в столбце 12, а остальные столбцы служат предикторами, полностью команда выглядит следующим образом:

```
> library(Cubist) > m.cubist <- cubist(x = wine_train[-12], y = wine_train$quality)
```

Для того чтобы проверить базовую информацию о дереве модели, нужно ввести его имя:

```
> m.cubistCall:cubist.default(x = wine_train[-12], y = wine_train$quality)
Number of samples: 3750
Number of predictors: 11
Number of committees: 1
Number of rules: 25
```

Как видим, алгоритм сгенерировал 25 правил для моделирования качества вина. Для того чтобы изучить некоторые из них, мы можем применить к объекту модели функцию `summary()`. Поскольку полное дерево слишком велико, здесь приведены только первые несколько строк вывода этой функции, описывающих первое правило принятия решения:

```
> summary(m.cubist)
Rule 1: [21 cases, mean 5.0, range 4 to 6, est err 0.5]
if free.sulfur.dioxide >
```

```

30      total.sulfur.dioxide >
195      total.sulfur.dioxide <=
235      sulphates > 0.64      alcohol >
9.1then      outcome = 573.6 + 0.0478
total.sulfur.dioxide      - 573 density
- 0.788 alcohol      + 0.186
residual.sugar - 4.73 volatile.acidity

```

Вероятно, вы заметили, что этот результат частично похож на регрессионное дерево, построенное нами ранее. Последовательность решений, принятых на основе винных свойств диоксида серы, сульфатов и алкоголя, создает правило, кульминацией которого и является окончательный прогноз. Однако ключевое отличие результатов этого дерева модели от предыдущих результатов, полученных с помощью регрессионного дерева, заключается в том, что здесь узлы заканчиваются не числовым прогнозом, а линейной моделью.

Линейная модель для этого правила представлена в выходных данных после ключевого слова `then` в виде оператора `outcome=`. Числа могут интерпретироваться так же, как модели множественной регрессии, построенные нами ранее. Каждое значение — это предполагаемая бета-характеристика соответствующего признака, то есть общее влияние этого признака на прогнозируемое качество вина. Например, коэффициент 0,186 для остаточного сахара подразумевает, что при увеличении значения остаточного сахара на 1 прогнозируемое качество вина увеличится на 0,186.

Важно отметить, что регрессионные эффекты, оцениваемые этой моделью, касаются только тех образцов вина, которые достигли этого узла; исследование всех результатов работы алгоритма Cubist показывает, что в этом дереве моделей построено в общей сложности 25 линейных моделей, по одной на каждое правило принятия решения; каждая из этих моделей имеет свои оценки параметров воздействия остаточного сахара, а также десяти других признаков.

Для того чтобы проверить эффективность этой модели, посмотрим, насколько хорошо она работает с неизвестными ей до сих пор тестовыми данными. Функция `predict()` возвращает вектор прогнозируемых значений:

```
> p.cubist <- predict(m.cubist, wine_test)
```

По-видимому, дерево моделей прогнозирует более широкий диапазон значений, чем регрессионное дерево:

```
> summary(p.cubist)  Min.  1st
Qu.   Median   Mean  3rd
```

Qu.	Max. 3.677	5.416	5.906	5.848	6.2
38	7.393				

Корреляция также, вероятно, существенно выше:

```
> cor(p.cubist, wine_test$quality)[1] 0.6201015
```

Более того, у этой модели средняя абсолютная ошибка немного меньше:

```
> MAE(wine_test$quality, p.cubist)[1] 0.5339725
```

Мы не достигли значительных улучшений по сравнению с регрессионным деревом, однако превысили эффективность модели нейронной сети, опубликованной Кортесом, и приблизились к опубликованному среднему значению абсолютной ошибки 0,45 для модели опорных векторов — и все это при использовании гораздо более простого метода обучения.



Мы подтвердили, что прогнозирование качества вин является сложной проблемой — что неудивительно. В конце концов, дегустация вина — дело субъективное. Если хотите дополнительно попрактиковаться, можете попытаться вернуться к этой задаче после того, как прочитаете главу 11, где рассматриваются дополнительные методы, которые могут привести к лучшим результатам.

Резюме

В этой главе мы рассмотрели два метода моделирования числовых данных. Первый из них, линейная регрессия, предполагает аппроксимацию данных прямой линией. Вторым методом используются деревья решений для числового прогнозирования. Эти деревья имеют две формы: регрессионные деревья, в которых используются средние значения примеров в концевых узлах для создания числовых прогнозов; и деревья моделей, в которых применяется гибридный подход с построением регрессионной модели для каждого концевого узла. Такой подход в некоторых отношениях работает лучше, чем каждый из использованных в нем методов по отдельности.

Мы начали понимать полезность регрессионного моделирования, применив его для расследования причин катастрофы космического челнока «Челленджер». Затем мы использовали линейное регрессионное моделирование для расчета ожидаемых медицинских расходов для различных слоев населения. Поскольку оценочная регрессионная модель хорошо описывает взаимосвязь между признаками и целевой переменной, нам удалось выявить некоторые группы людей, имеющих общие признаки. К ним относятся курильщики и тучные люди, которым, возможно, потребуется выставлять более высокие страховые тарифы, чтобы покрыть их медицинские расходы, превышающие средний показатель.

Регрессионные деревья и деревья моделей были использованы для моделирования субъективного представления о качестве вин на основе измеряемых характеристик. При этом мы узнали, что регрессионные деревья предлагают простой способ объяснить связь между признаками и числовым результатом, но для обеспечения высокой точности могут потребоваться более сложные деревья. Попутно мы изучили новые методы оценки эффективности числовых моделей.

Если в этой главе описаны методы машинного обучения, которые дают четкое понимание взаимосвязей между входными и выходными данными, то в следующей главе будут рассмотрены методы, которые приводят к практически непонятным моделям. Преимуществом этих методов является то, что они чрезвычайно мощные — относятся к числу самых мощных стандартных классификаторов — и могут применяться как к задачам классификации, так и к задачам числового прогнозирования.

7. Методы «черного ящика»: нейронные сети и метод опорных векторов

Выдающийся автор произведений научной фантастики Артур Кларк (Arthur C. Clarke) писал: «Любая достаточно развитая технология неотличима от магии». В этой главе рассмотрено несколько методов машинного обучения, которые на первый взгляд могут показаться волшебством. Они чрезвычайно мощны, однако их внутренний принцип работы трудно понять.

В инженерии такие методы называют процессами «*черного ящика*», так как механизм, который преобразует входные данные в выходные, скрыт внутри воображаемого блока. Так, программное обеспечение с закрытым исходным кодом преднамеренно скрыто внутри «*черного ящика*», в котором находятся защищенные патентами алгоритмы; в политическом законодательстве «*черный ящик*» — это бюрократические процессы, а при изготовлении колбасы «*черный ящик*» означает некоторое сознательное невежество относительно деталей приготовления этого продукта. При машинном обучении «*черный ящик*» обусловлен сложной математикой, обеспечивающей функционирование этих алгоритмов.

Понять, как работают эти методы, нелегко, однако слепо применять модели «*черного ящика*» опасно. В этой главе мы заглянем в середину «*ящика*» и исследуем весь статистический «*фарш*», используемый при построении таких моделей. Мы узнаем ответы на следующие вопросы.

- Как нейронные сети имитируют живой мозг для моделирования математических функций?
- Как в методе опорных векторов используются многомерные поверхности для определения взаимосвязи между признаками и результатами?

- Как, несмотря на всю сложность, эти методы можно легко применять для решения реальных задач?

Если вам повезет, то вы поймете, что для применения методов машинного обучения типа «черный ящик» не требуется черный пояс по статистике — просто не нужно бояться!

Нейронные сети

Искусственная нейронная сеть, ИНС (Artificial Neural Network, ANN) моделирует взаимосвязи между множеством входных сигналов и выходным сигналом, используя модель, основанную на принципе реагирования биологического мозга на раздражители, поступающие от органов чувств. Подобно тому как мозг использует сеть взаимосвязанных клеток — *нейронов* — для обеспечения широких возможностей обучения, *ИНС* задействует сеть искусственных нейронов, или *узлов*, для решения сложных задач обучения.

Человеческий мозг состоит примерно из 85 миллиардов нейронов, которые образуют сеть, способную интерпретировать невероятный объем знаний. Как и следовало ожидать, это многократно превышает возможности мозга остальных живых существ. Например, у кошки около миллиарда нейронов, у мыши — примерно 75 миллионов, а у таракана — всего лишь миллион. Что же касается большинства нейронных сетей, то у них гораздо меньше нейронов — как правило, всего несколько сотен. Именно поэтому мы не рискуем обещать появление искусственного мозга в ближайшем будущем: даже плодовая муха с ее 100 000 нейронов намного превосходит современную нейронную сеть.

Несмотря на то что современная нейронная сеть не способна полностью смоделировать даже мозг таракана, она может воспроизвести адекватную эвристическую модель его поведения. Предположим, что мы разрабатываем алгоритм, который бы имитировал, как убегает таракан, когда его обнаруживают. Если поведение робота-таракана будет убедительным, так ли уж важно, что его мозг уступает по сложности живому существу? На этом вопросе основан спорный *тест Тьюринга*, предложенный в 1950 году Аланом Тьюрингом (Alan Turing), пионером в области информатики. Тест считает машину мыслящей, если человек не в состоянии отличить ее поведение от поведения живого существа.



Подробнее об интригах и спорах вокруг теста Тьюринга читайте в Стэнфордской философской энциклопедии (<https://plato.stanford.edu/entries/turing-test/>).

Примитивные нейронные сети используются для симуляции решения задач мозгом вот уже более 50 лет. Поначалу при этом применялись такие

простые функции, как логическое «И» либо логическое «ИЛИ». Целью первых исследований было, прежде всего, помочь ученым понять, как работает биологический мозг. Однако, поскольку за последние годы компьютеры стали значительно мощнее, сложность нейронных сетей тоже возросла настолько, что теперь нейронные сети часто применяются для решения более практических задач, включая следующие:

- создание программ распознавания речи, рукописного текста и программ распознавания изображений, используемых, например, в приложениях для смартфонов, программ сортировки почты и поисковых систем;
- автоматизацию интеллектуальных устройств, таких как системы контроля климата в офисных зданиях или автопилоты в автомобилях и беспилотных летательных аппаратах;
- построение сложных погодных и климатических моделей, моделирование предела прочности, гидродинамики и многих других научных, социальных или экономических явлений.

Вообще, нейронные сети являются универсальными алгоритмами обучения, которые можно применять для решения практически любых задач обучения: классификации, числового прогнозирования и даже для распознавания образов без учителя.



Независимо от того, заслуживают они этого или нет, в средствах массовой информации об обучающих алгоритмах нейронных сетей часто сообщают с большой помпезностью. Например, «искусственный мозг», разработанный Google, хвалили за то, что он способен распознавать видео с кошками на YouTube. Такая шумиха обычно связана не столько с тем, что алгоритмы нейронных сетей представляют собой нечто уникальное, сколько с тем фактом, что нейронные сети являются неким подобием живого разума.

ИНС часто применяются для решения задач, в которых входные и выходные данные определены четко, однако процессы, их связывающие, чрезвычайно сложны и трудно поддаются определению. Нейронные сети, как и метод «черного ящика», хорошо работают для подобных задач.

От биологических нейронов — к искусственным

Поскольку нейронные сети специально разрабатывались как концептуальные модели деятельности человеческого мозга, сначала надо понять, как функционируют биологические нейроны. Как показано на рис. 7.1, входные сигналы принимаются *дендритами* клетки в результате биохимического процесса. Этот процесс позволяет определить силу (вес) импульса в соответствии с его относительной важностью или частотой. *Тело клетки* начинает накапливать поступающие сигналы, и в

какой-то момент достигается порог, при котором клетка срабатывает, выходной сигнал в результате электрохимического процесса передается дальше по *аксону*. На терминалях (концевых участках) аксона электрический сигнал снова преобразуется в химический, и этот химический сигнал передается соседним нейронам через крошечный зазор, называемый *синапсом*.

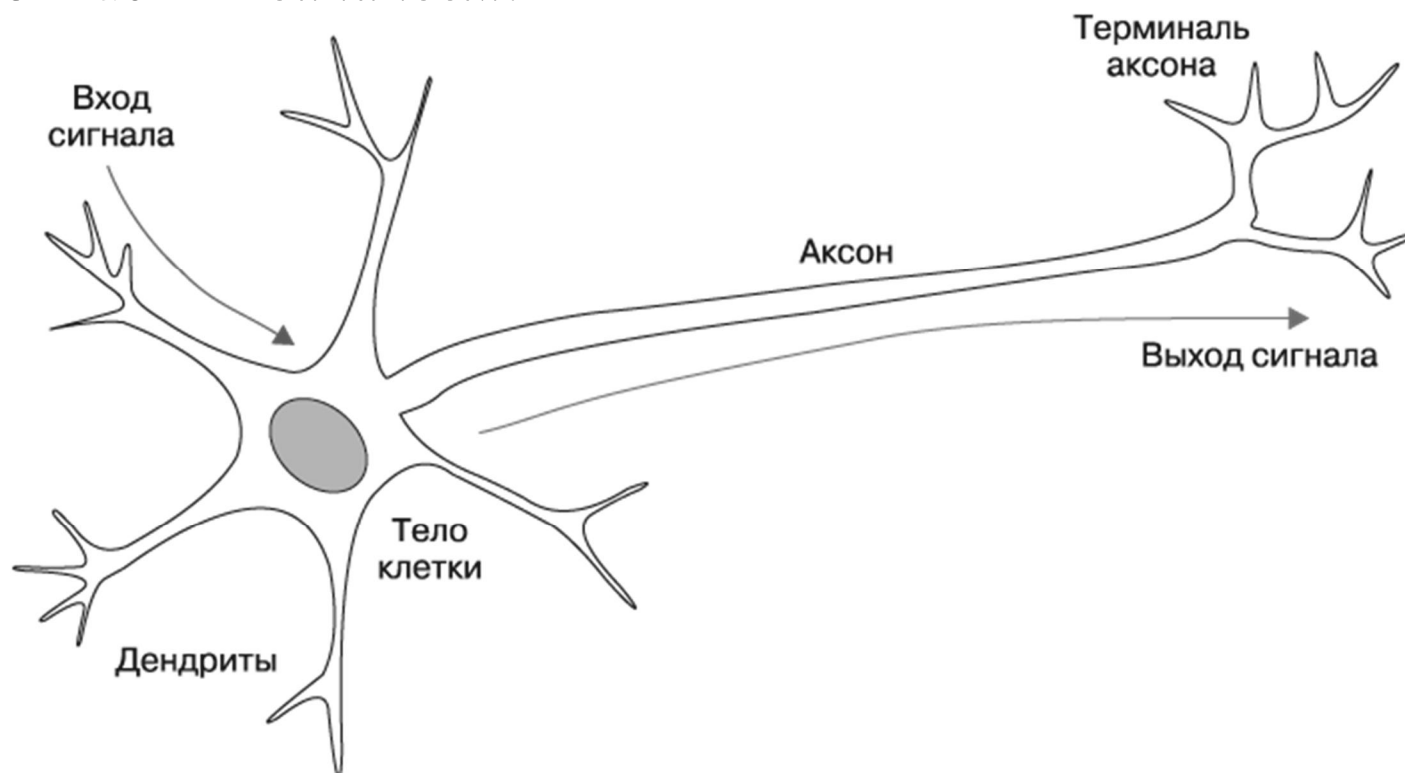


Рис. 7.1. Наглядное изображение биологического нейрона

Модель отдельного искусственного нейрона может быть описана в терминах, очень похожих на те, что используются в биологической модели. Как показано на рис. 7.2, направленная сетевая диаграмма определяет взаимосвязь между входными сигналами, полученными дендритами (переменными x), и выходным сигналом (переменной y). Как и в случае с биологическим нейроном, сигналу каждого дендрита присваивается вес (значения w) в соответствии с важностью — пока не будем обращать внимание на то, как именно определяются эти веса. В теле ячейки входные сигналы суммируются, и сигнал передается в соответствии с *функцией активации*, обозначенной буквой f .

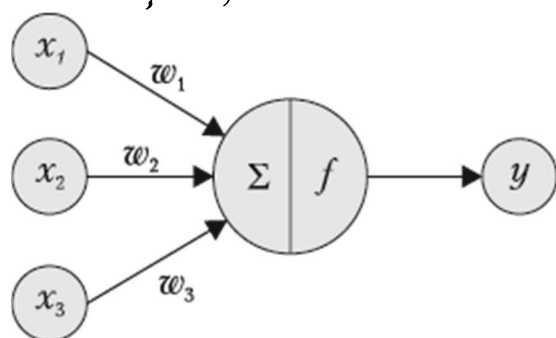


Рис. 7.2. Искусственный нейрон имитирует структуру и функции биологического нейрона

Типичный искусственный нейрон, имеющий n входных дендритов, может быть представлен формулой, приведенной далее. Весовые коэффициенты w позволяют каждому из n входов (обозначенных как x_i) вносить больший или меньший вклад в сумму входных сигналов. Полученная в итоге сумма используется функцией активации $f(x)$, а результирующий сигнал $y(x)$ является выходным аксоном:

$$y(x) = f\left(\sum_{i=1}^n w_i x_i\right).$$

В нейронных сетях нейроны, определенные таким образом, служат строительными блоками для создания сложных моделей данных.

Существует множество вариантов нейронных сетей, но все из них имеют такие характеристики, как:

- *функция активации*, которая преобразует суммарный входной сигнал нейрона в один выходной сигнал для дальнейшей передачи по сети;
- *топология сети* (архитектура), которая описывает количество нейронов модели, а также количество слоев и способ их соединения;
- *алгоритм обучения*, который определяет, как устанавливаются веса связей, чтобы тормозить или возбуждать нейроны пропорционально входному сигналу.

Рассмотрим несколько вариаций каждой из этих категорий, чтобы увидеть, как можно их использовать для построения типичных моделей нейронных сетей.

Функции активации

Функция активации — это механизм, с помощью которого искусственный нейрон обрабатывает входящую информацию и передает ее по сети. Так же как искусственный нейрон является подобием его биологической версии, так и функция активации моделируется аналогично природной.

Биологическую функцию активации можно представить как процесс, который включает в себя суммирование общего входного сигнала и определение того, проходит ли он порог срабатывания. Если проходит, то нейрон передает сигнал; в противном случае ничего не происходит. В терминологии нейронных сетей это называется *пороговой функцией активации*, поскольку она выдает выходной сигнал только при условии превышения входными сигналами заданного порога.

На рис. 7.3 изображена типичная пороговая функция; в данном случае нейрон срабатывает тогда, когда сумма входных сигналов не меньше нуля. Поскольку форма функции напоминает ступеньку, ее иногда называют функцией активации единичного скачка.

Пороговая функция активации интересна своим подобием биологическому аналогу, однако в нейронных сетях она используется редко. Функции активации нейронных сетей, свободные от ограничений биохимии, могут строиться на основе их способности демонстрировать желаемые математические характеристики и точно моделировать отношения между данными.

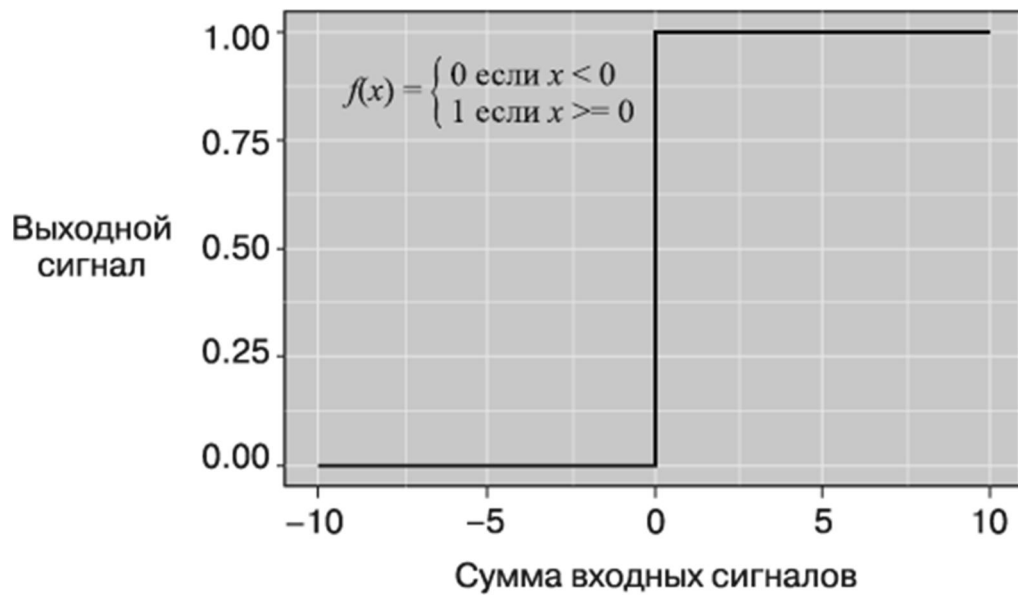


Рис. 7.3. Пороговая функция активации активируется только тогда, когда входные сигналы превышают пороговое значение

Возможно, наиболее распространенная альтернатива — это *сигмоидная функция активации* (а именно, *логистическая сигмоида*), показанная на рис. 7.4. Обратите внимание, что в представленной формуле e — это основание натурального логарифма (примерно равное 2,72). Несмотря на то что эта функция имеет похожую ступенчатую, S-образную, форму, как и пороговая функция активации, выходной сигнал больше не является двоичным; выходные значения могут принимать любые значения в диапазоне от 0 до 1.

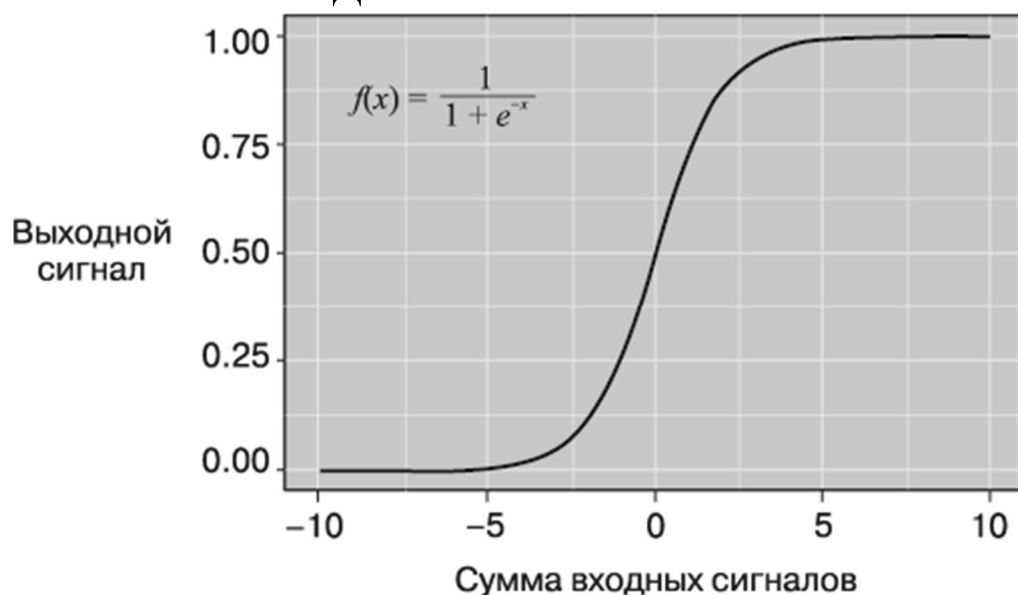


Рис. 7.4. Сигмоидная функция активации имитирует биологическую функцию активации в виде гладкой кривой

Кроме того, сигмоидная функция является *дифференцируемой*, а это означает, что можно вычислить производную по всему диапазону входных данных. Как вы скоро убедитесь, это свойство имеет решающее значение для создания эффективных алгоритмов оптимизации нейронных сетей.

Несмотря на то что сигмоидная функция является, пожалуй, наиболее используемой функцией активации и часто применяется по умолчанию, некоторые алгоритмы нейронных сетей допускают выбор альтернатив. Несколько таких функций активации показаны на рис. 7.5.

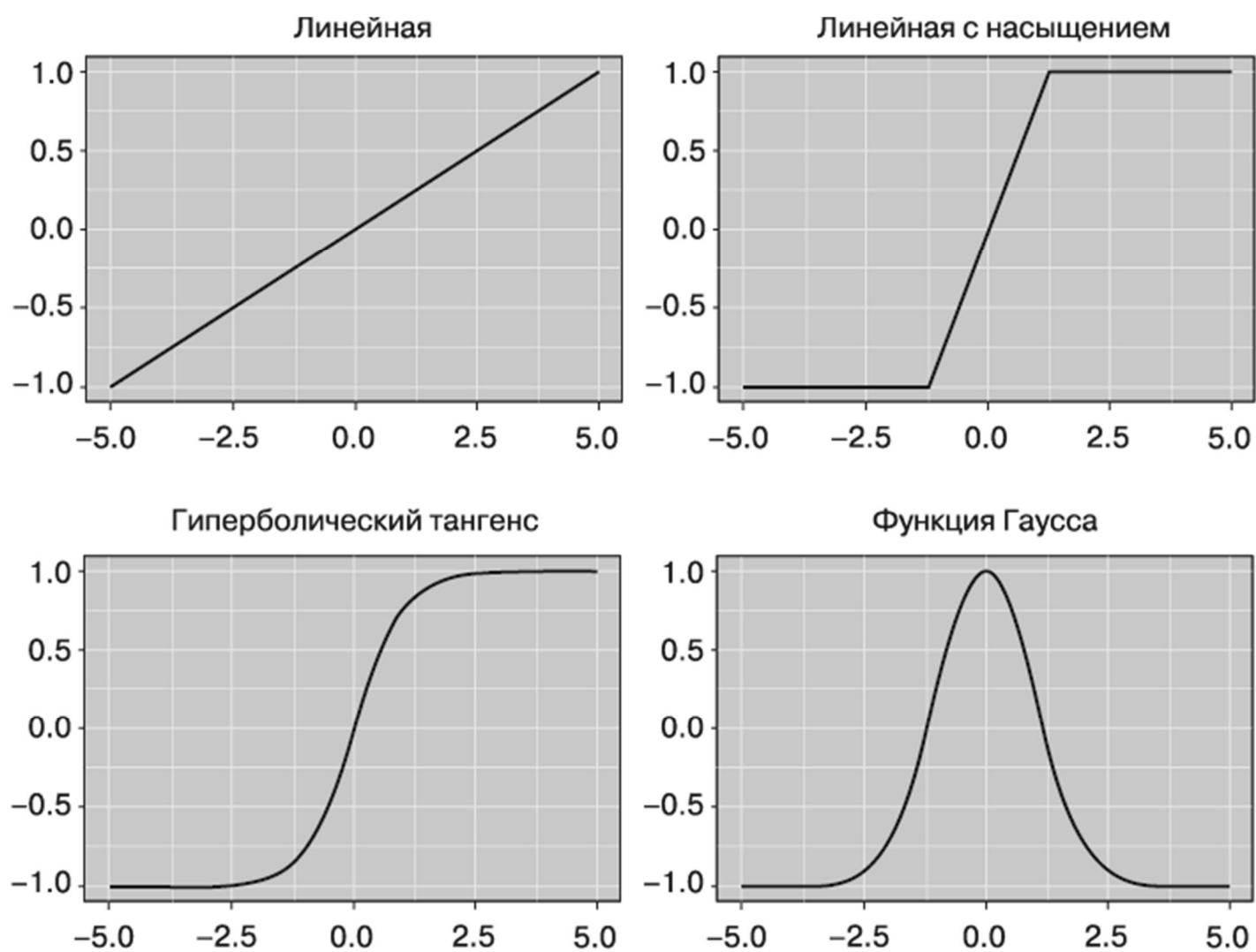


Рис. 7.5. Некоторые распространенные функции активации нейронных сетей

Главная отличительная деталь этих функций активации — диапазон значений выходного сигнала. Обычно он составляет $(0, 1)$, $(-1, +1)$ или $(-\infty, +\infty)$. Выбор функции активации влияет на нейронную сеть таким образом, чтобы она наиболее точно соответствовала определенным типам данных. Это позволяет создавать специализированные нейронные сети.

Например, линейная функция активации приводит к нейронной сети, очень похожей на модель линейной регрессии, а гауссова функция активации является основой для сетей *радиально-базисных функций* (Radial Basis Function, RBF). У каждой из этих функций есть свои сильные стороны, которые для одних задач обучения подходят лучше, а для других — хуже.

Важно понимать, что у многих функций активации диапазон входных значений, которые влияют на выходной сигнал, относительно узок. Например, в случае сигмоидной функции для всех значений входного сигнала ниже -5 выходной сигнал очень близок к 0, а для всех значений входного сигнала выше $+5$ — очень близок к 1. Такое сжатие сигнала приводит к насыщению выходного сигнала для верхних и нижних значений диапазона очень динамичных входных сигналов — точно так же, как повышение громкости гитарного усилителя приводит к искажению звука из-за обрезания пиков звуковых волн. Поскольку это, по сути, сжимает входные значения в меньший диапазон выходных данных, функции активации, подобные сигмоидной, иногда называют *функциями сжатия*.

Одно из решений проблемы сжатия — преобразование входных сигналов нейронной сети таким образом, чтобы значения всех признаков

находились в небольшом диапазоне около 0. Решение может включать стандартизацию или нормализацию признаков. При ограничении диапазона входных значений функция активации будет действовать для всего диапазона. Преимуществом является то, что модель также может быстрее обучаться, поскольку алгоритм будет быстрее выполнять итерацию для рабочего диапазона входных значений.



Теоретически нейронная сеть способна адаптироваться к очень динамично изменяющимся признакам, регулируя их вес на протяжении многих итераций. Однако многие алгоритмы прекратят итерации задолго до того, как это произойдет. Если ваша модель не сходится, проверьте еще раз, правильно ли вы стандартизировали входные данные. Также, возможно, стоит выбрать другую функцию активации.

Топология сети

Способность нейронной сети к обучению основана на ее *топологии* — паттернах и структурах взаимосвязанных нейронов. Существует множество форм сетевой архитектуры, однако они различаются по таким трем ключевым характеристикам, как:

- количество слоев;
- указание, может ли информация перемещаться по сети в обратном направлении;
- количество узлов в каждом слое сети.

Топология определяет сложность задач, которые могут быть решены с помощью сети. Как правило, чем больше и сложнее сеть, тем более неявные закономерности она может находить и тем более широкие границы принятия решений имеет. Однако возможности сети зависят не только от ее размера, но и от способа расположения узлов.

Количество слоев

Чтобы описать топологию, нам понадобится терминология, которая позволила бы различать искусственные нейроны в зависимости от их расположения в сети. На рис. 7.6 показана топология очень простой сети. Множество нейронов, называемых *входными узлами*, получает необработанные сигналы непосредственно из входных данных. Каждый входной узел отвечает за обработку одного признака из набора данных; затем значение признака преобразуется функцией активации соответствующего узла. Сигналы, отправляемые входными узлами, принимаются *выходным узлом*, который использует собственную функцию активации для генерации окончательного прогноза (обозначена здесь буквой *p*).

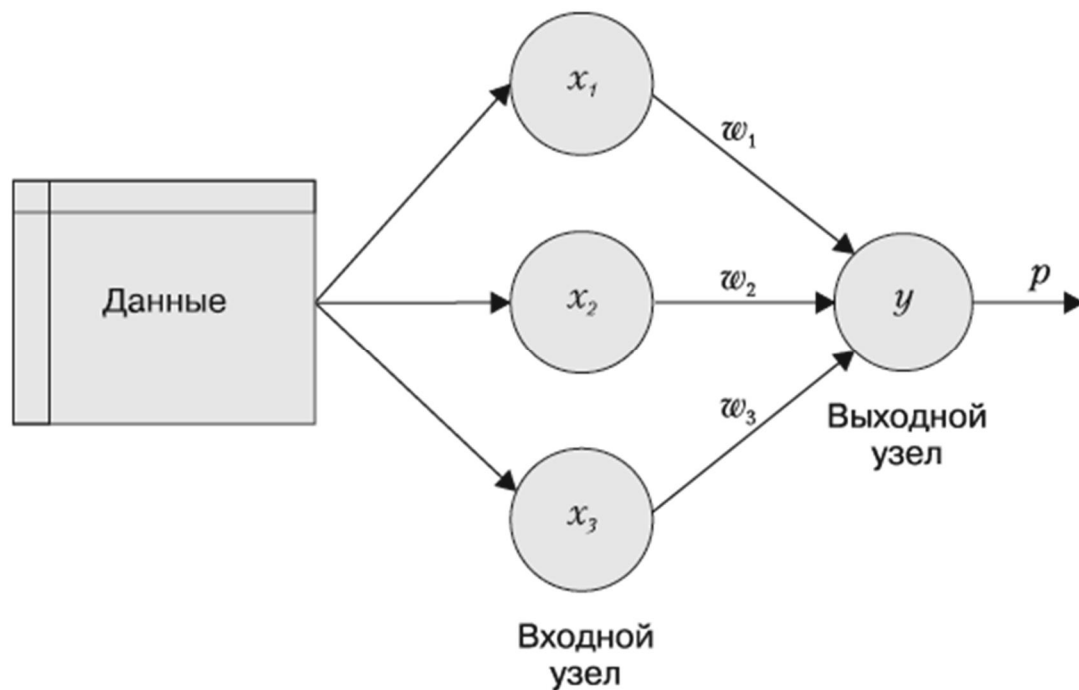


Рис. 7.6. Простая однослойная нейронная сеть с тремя входными узлами

Входные и выходные узлы объединены в группы — *слои*. Поскольку входные узлы обрабатывают входные данные точно в таком виде, в каком они поступают, сеть имеет только один набор весов связей (обозначены здесь как w_1 , w_2 и w_3), поэтому такая сеть называется *однослойной*. Однослойные сети можно использовать для базовой классификации образов, особенно тех, которые линейно разделимы, но для большинства задач обучения нужны более сложные сети.

Нетрудно догадаться, что явный способ создания более сложных сетей — введение дополнительных слоев. Как показано на рис. 7.7, *многослойная сеть* — это сеть, в которую добавлены один или несколько *скрытых слоев*. Они обрабатывают сигналы, поступающие от входных узлов, пока те не достигнут выходного узла. Большинство многослойных сетей являются полносвязными. Это означает, что каждый узел одного уровня связан с каждым узлом следующего уровня, но это опционально.

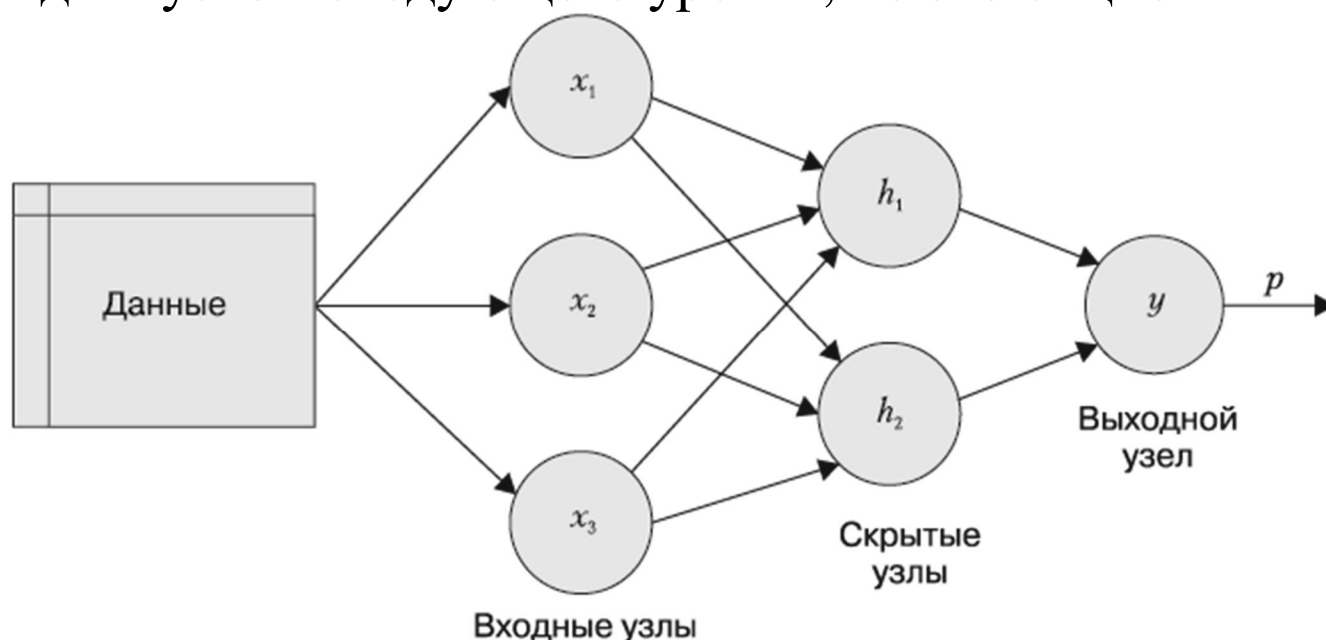


Рис. 7.7. Многослойная сеть с одним скрытым слоем, состоящим из двух узлов

Направление перемещения информации

Возможно, вы обратили внимание, что в предыдущих примерах стрелки обозначают сигналы, которые распространяются только в одном направлении. Сети, в которых входной сигнал непрерывно передается в

одном направлении от входного уровня к выходному, называются *сетями прямого распространения сигнала* (feedforward networks).

Несмотря на ограничение передачи информации, сети прямого распространения обеспечивают удивительную гибкость. Например, можно варьировать количество уровней и узлов на каждом уровне, моделировать несколько выходных сигналов одновременно или использовать несколько скрытых слоев (рис. 7.8). Нейронная сеть с несколькими скрытыми слоями называется *глубокой нейронной сетью* (Deep Neural Network, DNN), а практика обучения таких сетей — *глубоким обучением*. Глубокие нейронные сети, обученные на больших наборах данных, способны при решении сложных задач, таких как распознавание изображений и обработка текста, действовать подобно человеку.

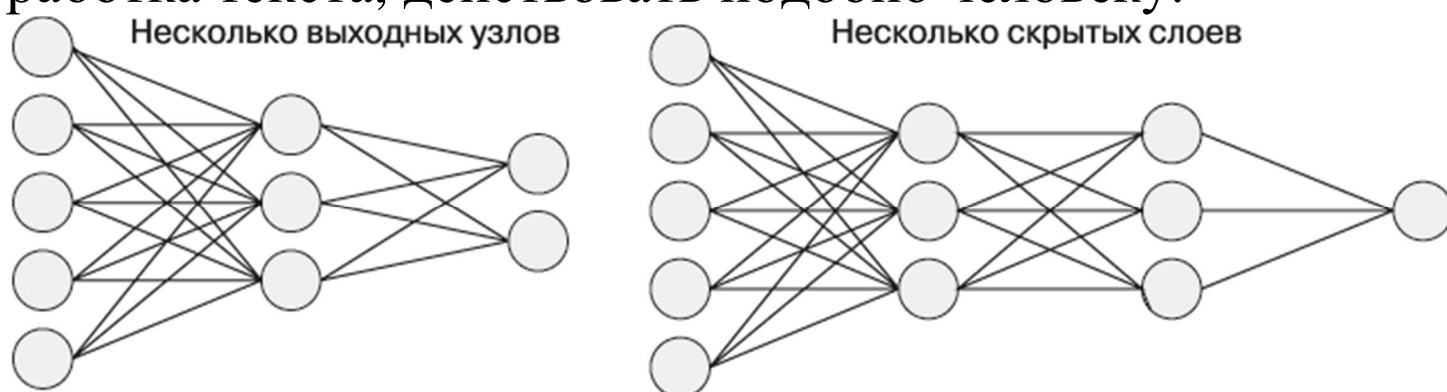


Рис. 7.8. Сложные нейронные сети могут иметь несколько выходных узлов и несколько скрытых слоев

В отличие от сетей прямого распространения, *рекуррентные сети* (или *сети с обратной связью*) позволяют передавать сигналы назад, образуя петли. Это свойство более точно отражает работу биологической нейронной сети и позволяет изучать чрезвычайно сложные образы. Добавление кратковременной памяти, или *задержки*, еще больше увеличивает возможности рекуррентных сетей. В частности, появляется возможность распознавать последовательность событий, происходящих в течение некоего отрезка времени. Это может использоваться для прогнозирования событий на фондовом рынке, распознавания речи или составления прогноза погоды. На рис. 7.9 показано, как выглядит простая рекуррентная сеть.

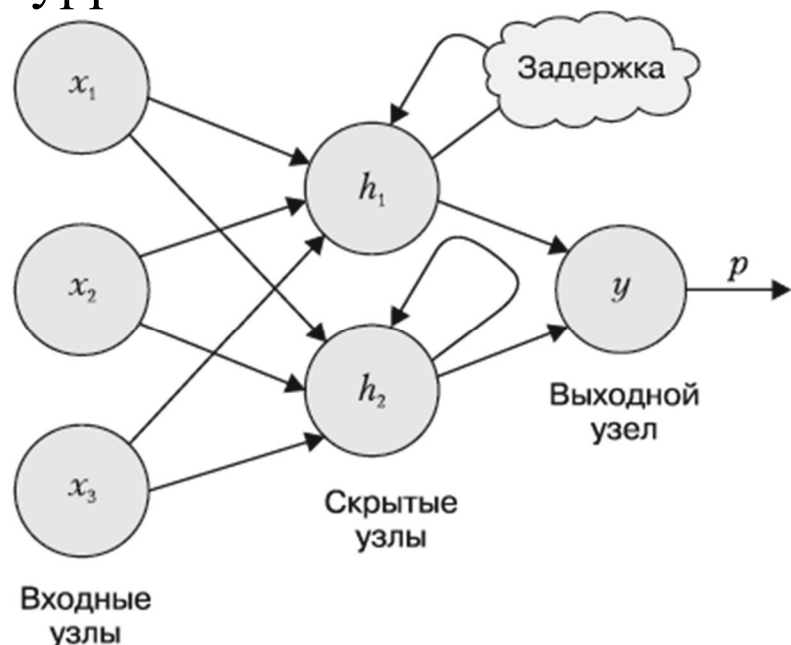


Рис. 7.9. Если допустить распространение информации назад по сети, то можно моделировать временную задержку

DNN и рекуррентные сети все чаще применяются в различных важных приложениях и, следовательно, становятся очень популярными. Однако при построении таких сетей используются методы и программное обеспечение, изучение которых выходит за рамки данной книги; кроме того, такие сети часто требуют доступа к специализированному компьютерному оборудованию или облачным серверам. Впрочем, более простые сети прямого распространения также способны моделировать многие реальные задачи. В сущности, многослойная сеть прямого распространения, или *многослойный перцептрон* (Multilayer Perceptron, MLP), является де-факто стандартной топологией нейронных сетей. Если вы заинтересовались глубоким обучением, то понимание топологии MLP обеспечит вам сильную теоретическую базу для построения более сложных моделей DNN в будущем.

Количество узлов в слое

Помимо разницы в количестве слоев и направлении передачи информации, нейронные сети различаются по сложности — по количеству узлов в каждом слое. Количество входных узлов определяется количеством признаков во входных данных. Аналогично количество выходных узлов определяется числом моделируемых результатов или количеством уровней классов в результате. Однако решение о количестве скрытых узлов принимает пользователь, прежде чем приступить к обучению модели.

К сожалению, не существует надежного правила определения количества нейронов в скрытом слое. Помимо многих других факторов, число таких нейронов зависит от количества входных узлов, объема тренировочных данных, количества зашумленных данных и сложности задачи обучения.

В целом чем сложнее топология и чем больше в ней сетевых связей, тем более сложные задачи обучения она позволяет решать. Больше количество нейронов приводит к построению модели, которая более точно отражает тренировочные данные, однако при этом существует риск переобучения; модель может плохо обобщать данные в будущем. Большие нейронные сети также могут оказаться чересчур затратными в вычислительном отношении и медленными в обучении.

Рекомендуется использовать минимальное количество узлов, которое обеспечивает достаточную эффективность для тестового набора данных. В большинстве случаев даже с небольшим количеством скрытых узлов — зачастую всего несколько штук — удастся построить нейронную сеть с впечатляющими способностями к обучению.



Доказано, что нейронная сеть, которая имеет по крайней мере один скрытый слой, насчитывающий достаточное количество нейронов, является универсальным аппроксиматором функции. Это означает, что нейронные сети могут использоваться для аппроксимации любой непрерывной функции с произвольной точностью на конечном интервале.

Обучение нейронной сети методом обратного распространения ошибки

Топология сети — это чистый лист, который сам по себе еще ничему не обучился. Он, как и новорожденный ребенок, учится на собственном опыте. Поскольку нейронная сеть обрабатывает входные данные, связи между нейронами усиливаются или ослабевают подобно тому, как развивается мозг ребенка, когда он знакомится с окружающей средой. Веса сетевых связей корректируются, отражая закономерности, наблюдаемые в течение определенного времени.

Обучение нейронной сети путем корректировки весов связей требует больших вычислительных ресурсов. Из-за этого нейронные сети, хотя они и были изучены несколько десятилетий назад, редко применялись для решения реальных задач до второй половины 1980-х годов, когда был изобретен эффективный метод обучения нейронных сетей. Этот алгоритм, в котором была использована стратегия обратного распространения ошибок, сейчас называется *методом обратного распространения ошибки* (backpropagation).



Случилось так, что метод обратного распространения ошибки был изобретен и опубликован несколькими исследовательскими группами примерно в одно и то же время, независимо друг от друга. Пожалуй, наиболее часто цитируемая из этих работ — Rumelhart D.E., Hinton G.E., Williams R.J. Learning representations by back-propagating errors. *Nature*, 1986. Vol. 323. P. 533–566.

Несмотря на то что метод обратного распространения ошибки все еще остается дорогостоящим по сравнению со многими другими алгоритмами ML, он поспособствовал возобновлению интереса к нейронным сетям. В результате в области интеллектуального анализа данных приобрели широкую популярность многослойные сети прямого распространения, в которых реализован метод обратного распространения ошибки. Его преимущества и недостатки перечислены в табл. 7.1.

Таблица 7.1

Преимущества	Недостатки
Алгоритм может быть адаптирован для решения задач классификации или числового прогнозирования. Позволяет моделировать более сложные паттерны, чем любой другой алгоритм. Делает несколько предположений об основных взаимосвязях между данными	Чрезвычайно сложный в вычислительном отношении и медленный процесс обучения, особенно при сложной топологии сети. Склонен к переобучению на тренировочных данных. В результате получается сложная модель типа «черный ящик», которую трудно, а иногда и невозможно интерпретировать

В наиболее обобщенной форме алгоритм обратного распространения ошибки проходит множество циклов, каждый из которых состоит из двух процессов. Один такой цикл называется *эпохой*. Поскольку сеть не содержит *априорных* (существующих) знаний, начальные веса обычно задаются случайным образом. Затем алгоритм в течение нескольких итераций выполняет процессы, пока не будет достигнут критерий остановки. Каждая эпоха алгоритма обратного распространения ошибки включает в себя:

- фазу прямого распространения сигнала, в которой нейроны активируются последовательно от входного слоя к выходному, по пути применяя к сигналам веса каждого нейрона и функцию активации. Когда сигналы достигают последнего слоя, выдается выходной сигнал;
- фазу обратного распространения ошибки, в которой выходной сигнал сети, полученный во время фазы прямого распространения, сравнивается с реальным целевым значением тренировочных данных. Разница между выходным сигналом сети и истинным значением дает ошибку, которая распространяется по сети в обратном направлении, чтобы изменить вес связей между нейронами и уменьшить будущие ошибки.

По мере выполнения алгоритм использует информацию, передаваемую в обратном направлении, чтобы уменьшить общую ошибку сети. И все же остается вопрос: поскольку взаимосвязи между входами и выходами нейронов так сложны, каким образом алгоритм определяет, какой именно вес следует изменить? Ответ на этот вопрос дает техника, называемая *градиентным спуском*. Концептуально эта техника работает аналогично тому, как путешественник, потерявшись в джунглях, находит путь к воде. Исследуя местность и непрерывно двигаясь в направлении с наибольшим спуском, он в итоге достигнет самой низкой впадины, которая, вероятно, окажется руслом реки.

У алгоритма обратного распространения ошибки есть аналогичный процесс, в ходе которого для определения градиента в направлении каждого из весов входящих сигналов используется производная функции активации нейрона — именно поэтому так важна дифференцируемая функция активации. Градиент показывает, насколько сильно будет уменьшаться или увеличиваться ошибка при изменении веса. Алгоритм будет пытаться изменять те веса, которые приводят к наибольшему

уменьшению ошибки, — эта величина называется *скоростью обучения*. Чем выше скорость обучения, тем быстрее алгоритм пытается спуститься по градиентам, что может сократить время обучения, однако повышается риск пропустить впадину (рис. 7.10).

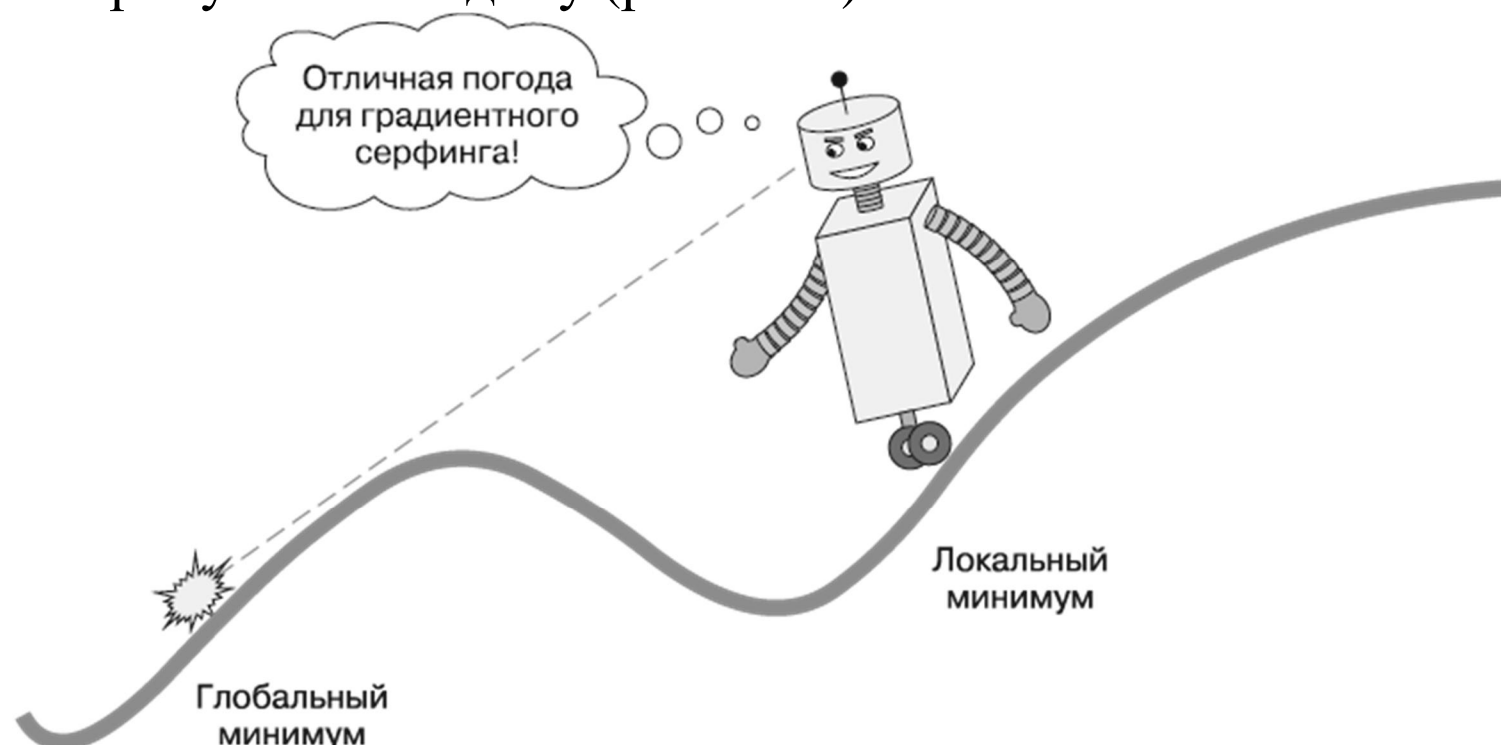


Рис. 7.10. Алгоритм градиентного спуска пытается найти минимальную ошибку, но может попасть на локальный минимум

Этот процесс представляется сложным, однако на практике он легко реализуется. Применим наши знания о многослойных сетях прямого распространения для решения реальной задачи.

Пример: моделирование прочности бетона с помощью нейронной сети

В области инженерно-строительных работ крайне важно иметь точные оценки эффективности строительных материалов. Эти оценки необходимы для разработки правил безопасности, регулирующих использование материалов при строительстве зданий, мостов и дорог.

Особый интерес представляет оценка прочности бетона. Бетон используется практически при любом строительстве, эксплуатационные характеристики бетона сильно различаются, так как он состоит из огромного количества ингредиентов, которые взаимодействуют комплексно. В результате трудно точно сказать, какова будет прочность готового продукта. Модель, которая бы позволяла определить прочность бетона наверняка, с учетом состава исходных материалов, могла бы обеспечить более высокий уровень безопасности строительных объектов.

Шаг 1. Сбор данных

Для этого анализа мы будем использовать данные о прочности бетона на сжатие, переданные Ай-Ченг Йе (I-Cheng Yeh) в репозиторий данных для машинного обучения *UCI Machine Learning Repository* (<http://archive.ics.uci.edu/ml>). Поскольку Ай-Ченг Йе успешно

использовал нейронные сети для моделирования этих данных, мы попытаемся воспроизвести его работу, применив простую модель нейронной сети в среде R.



Подробнее о подходе к этой задаче обучения Ай-Ченг Йе читайте в статье: Yeh I.C. Modeling of Strength of High-Performance Concrete Using Artificial Neural Networks. Cement and Concrete Research, 1998. Vol. 28. P. 1797–1808. Судя по сайту, этот набор данных содержит 1030 записей о разных марках бетона с восемью характеристиками, описывающими компоненты, используемые в составе бетонной смеси. Считается, что эти характеристики влияют на итоговую прочность при сжатии. К ним относятся: количество (в килограммах на кубический метр) цемента, воды, разных добавок, крупных и мелких заполнителей типа щебня и песка, используемых в готовом продукте, а также время схватывания (в днях).



Для того чтобы выполнить этот пример, загрузите файл `concrete.csv` и сохраните его в рабочем каталоге R.

Шаг 2. Исследование и подготовка данных

Как обычно, начнем анализ с загрузки данных в R-объект с помощью функции `read.csv()` и убедимся, что результат соответствует ожидаемой структуре:

```
> concrete <- read.csv("concrete.csv")>
str(concrete) 'data.frame':   1030 obs. of  9
variables: $ cement          : num 141 169 250 266 155
...$ slag                   : num 212 42.2 0 114 183.4 ...$
ash                          : num 0 124.3 95.7 0 0 ...$
water                        : num 204 158 187 228 193 ...$
superplastic: num 0 10.8 5.5 0 9.1 0 0 6.4 0 9
...$ coarseagg              : num 972 1081 957 932 1047 ...$
fineagg                      : num 748 796 861 670 697 ...$
age                           : int 28 14 28 28 28 90 7 56 28 28
...$ strength               : num 29.9 23.5 29.2 45.9 18.3
...
```

Девять переменных во фрейме данных соответствуют восьми характеристикам и одному ожидаемому результату, но стало очевидно, что есть проблема. Нейронные сети лучше всего работают тогда, когда входные данные масштабируются до узкого диапазона с центром около 0, а здесь мы видим значения в диапазоне от 0 до более чем 1000.

Как правило, решение такой проблемы заключается в изменении масштаба данных с помощью функции нормализации или стандартизации. Если распределение данных соответствует колоколообразной кривой (нормальное распределение, см. главу 2), то, возможно, имеет смысл использовать стандартизацию с помощью встроенной функции `scale()`. Если же распределение данных близко к равномерному или сильно отличается от нормального, то более подходящей может быть нормализация к диапазону от 0 до 1. В данном случае будем использовать последний вариант.

В главе 3 мы создали собственную функцию `normalize()`:

```
> normalize <- function(x) { return((x -  
min(x)) / (max(x) - min(x))) }
```

После того как будет выполнен этот код, можно применить функцию `normalize()` ко всем столбцам выбранного фрейма данных с помощью функции `lapply()`:

```
> concrete_norm <-  
as.data.frame(lapply(concrete, normalize))
```

Чтобы убедиться, что нормализация сработала, можно проверить, равны ли минимальное и максимальное значения признака `strength0` и `1` соответственно:

```
> summary(concrete_norm$strength)  Min.   1st  
Qu.   Median   Mean   3rd  
Qu.   Max. 0.0000  0.2664  0.4001  0.4172  
0.5457  1.0000
```

Для сравнения: исходные минимальное и максимальное значения этого признака были равны `2.33` и `82.60` соответственно:

```
> summary(concrete$strength)  Min.   1st  
Qu.   Median   Mean   3rd  
Qu.   Max.   2.33   23.71  34.44  35.82  
46.14  82.60
```



Любое преобразование, примененное к данным до обучения модели, впоследствии должно быть применено в обратном порядке, чтобы преобразовать признак назад в исходные единицы измерения. Для облегчения масштабирования целесообразно сохранить исходные данные или хотя бы сводную статистику исходных данных.

Следуя сценарию, описанному Йе в исходной статье, разделим данные на тренировочный набор, включающий в себя 75 % всех примеров, и тестовый набор, состоящий из 25 %. Используемый CSV-файл

отсортирован в случайном порядке, поэтому нам остается лишь разделить его на две части:

```
> concrete_train <- concrete_norm[1:773, ]>  
concrete_test <- concrete_norm[774:1030, ]
```

Мы будем использовать тренировочный набор данных для построения нейронной сети и тестовый набор данных для оценки того, насколько хорошо модель обобщается для будущих результатов. Поскольку нейронную сеть легко довести до состояния переобучения, этот шаг очень важен.

Шаг 3. Обучение модели на данных

Чтобы смоделировать взаимосвязь между ингредиентами, используемыми в производстве бетона, и прочностью готового продукта, построим многослойную нейронную сеть прямого распространения.

Пакет `neuralnet`, разработанный Стефаном Фритчем (Stefan Fritsch) и Фрауке Гюнтер (Frauke Guenther), обеспечивает стандартную и простую в использовании реализацию таких сетей. В этот пакет также входит функция для построения топологии сети. Реализация `neuralnet` — хороший способ получить дополнительную информацию о нейронных сетях, хотя это не означает, что ее нельзя использовать и для выполнения реальной работы — как вы скоро убедитесь, это довольно мощный инструмент.



В R есть еще несколько пакетов, обычно используемых для обучения моделей нейронных сетей, каждый из которых имеет свои достоинства и недостатки. Поскольку `nnet` поставляется в составе стандартного комплекта R, он является, пожалуй, наиболее часто упоминаемой реализацией нейронных сетей. В пакете использован немного более сложный алгоритм, чем стандартный метод обратного распространения ошибки. Еще один вариант — пакет `RSNNS`, содержащий полный набор функций для нейронной сети, но его недостатком является то, что этот пакет труднее освоить.

Поскольку пакет `neuralnet` не включен в базовый R, вам нужно будет его установить, введя

```
команду install.packages("neuralnet"), и загрузить с  
помощью команды library(neuralnet). Входящую в пакет  
функцию neuralnet() можно применять для обучения нейронных  
сетей числовому прогнозированию с использованием следующего  
синтаксиса.
```

Синтаксис нейронной сети

Использование функции `neuralnet()` из пакета `neuralnet`

Построение модели:

```
m <- neuralnet(target ~ predictors, data = mydata,  
              hidden = 1, act.fct = "logistic")
```

`target` – модель, которая будет построена в результате обучения на фрейме данных `mydata`;

`predictors` – R-формула, определяющая признаки из фрейма данных `mydata`, которые будут использоваться при прогнозе;

`data` – фрейм данных, которому принадлежат `target` и `predictors`;

`hidden` – количество нейронов в скрытом слое (по умолчанию 1). *Примечание:* для описания нескольких скрытых слоев можно использовать вектор целых чисел, например `c(2,2)`;

`act.fct` – функция активации: "logistic" или "tanh". *Примечание:* также может быть использована любая другая функция.

Функция возвращает объект нейронной сети, который можно использовать для прогнозирования.

Прогнозирование:

```
p <- compute(m, test)
```

`m` – модель, обученная с помощью функции `neuralnet()`;

`test` – фрейм данных, содержащий тестовые данные с теми же признаками, что и у тренировочных данных, и для построения классификатора.

Функция возвращает список, состоящий из двух компонентов: `$neurons`, где хранятся нейроны для каждого слоя сети, и `$net.result`, где хранятся значения, спрогнозированные с помощью данной модели.

Примеры:

```
concrete_model <- neuralnet(strength ~ cement + slag + ash,  
                           data = concrete, hidden = c(5, 5), act.fct = "tanh")  
model_results <- compute(concrete_model, concrete_data)  
strength_predictions <- model_results$net.result
```

Начнем с обучения простейшей многоуровневой сети прямого распространения с параметрами по умолчанию, имеющей только один скрытый узел:

```
> concrete_model <- neuralnet(strength ~ cement  
+ slag + ash + water + superplastic +  
coarseagg + fineagg + age, data =  
concrete_train)
```

Затем, как показано на рис. 7.11, можно визуализировать топологию сети, используя функцию `plot()` и передав ей полученный в результате объект модели:

```
> plot(concrete_model)
```

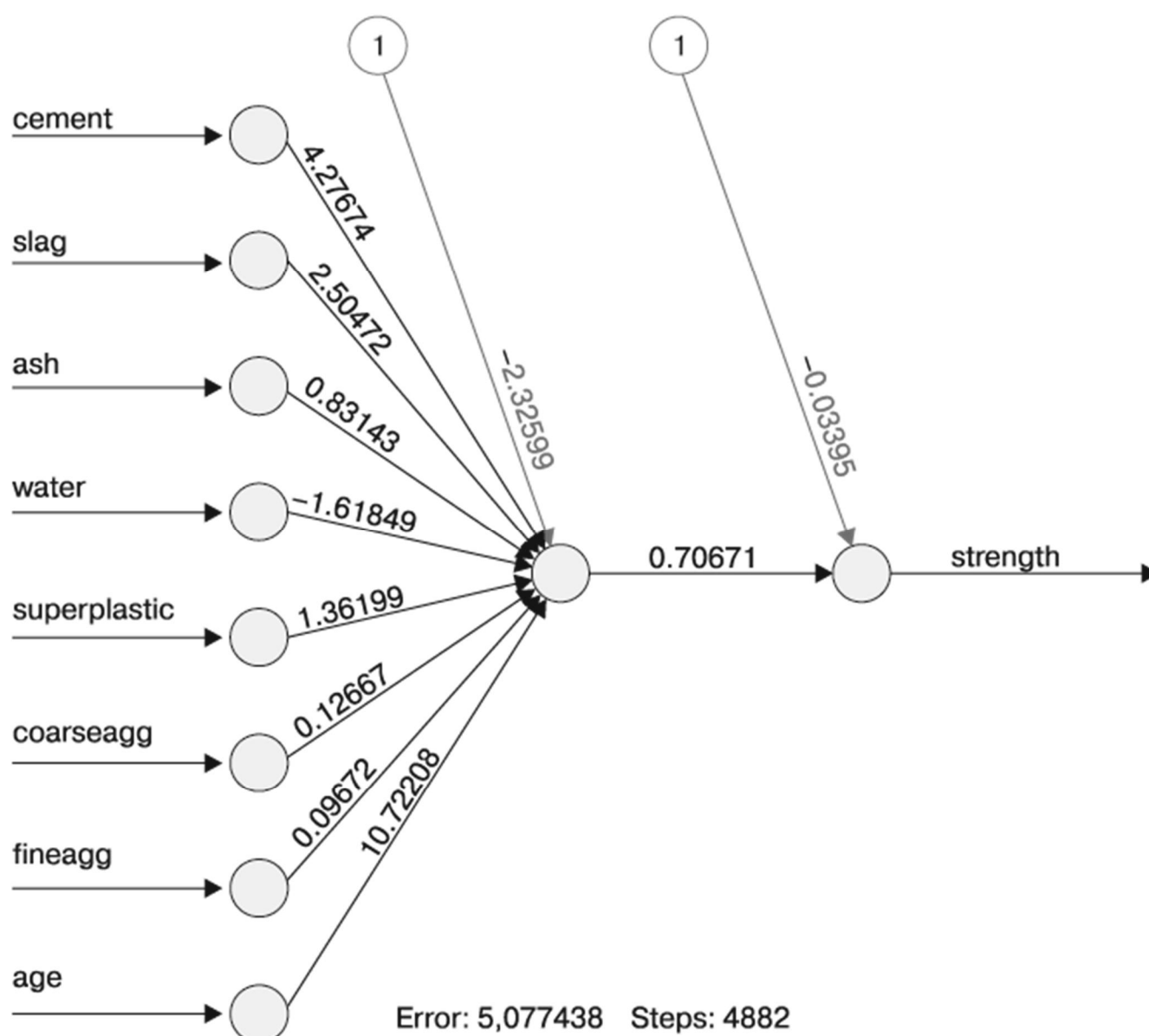



Рис. 7.11. Визуализация топологии нашей простой многоуровневой сети прямого распространения

В этой простой модели есть по одному входному узлу для каждого из восьми признаков, затем идет один скрытый и один выходной узел, который и дает прогноз прочности бетона. На диаграмме также изображены веса для каждой связи и *величина смещения*, обозначенная для узлов, помеченных цифрой 1. Величина смещения — это числовая константа, которая позволяет смещать значение в указанном узле вверх или вниз, примерно как сдвиг в линейном уравнении.



Нейронную сеть с одним скрытым узлом можно считать «двоюродной сестрой» моделей линейной регрессии, рассмотренных в главе 6. Веса между входными узлами и скрытым узлом подобны бета-коэффициентам, а вес смещения — сдвигу.

В нижней части рисунка отображаются количество шагов обучения и величина ошибки — *суммарная среднеквадратичная ошибка* (Sum of Squared Errors, SSE), которая, как и следовало ожидать, является суммой квадратов разностей между прогнозируемыми и фактическими значениями. Чем меньше SSE, тем точнее модель соответствует тренировочным данным, что свидетельствует об эффективности этих данных, но мало говорит о том, как модель будет работать с неизвестными данными.

Шаг 4. Оценка эффективности модели

Диаграмма топологии сети дает возможность заглянуть в «черный ящик» нейронной сети, но дает не очень много информации о том, насколько хорошо модель соответствует будущим данным. Для генерации прогнозов на тестовом наборе данных можно воспользоваться функцией `compute()`:

```
> model_results <- compute(concrete_model,
concrete_test[1:8])
```

Функция `compute()` работает немного иначе, чем использованные нами до сих пор функции `predict()`. Она возвращает список, состоящий из двух компонентов: `$neurons`, где хранятся нейроны для каждого слоя сети, и `$net.result`, где хранятся спрогнозированные значения. Именно `$net.result` нам и нужен:

```
> predicted_strength <-
model_results$net.result
```

Поскольку у нас задача числового прогнозирования, а не классификации, мы не можем использовать для проверки точности модели матрицу несоответствий. Измерим корреляцию между прогнозируемым и истинным значением прочности бетона. Если прогнозируемые и фактические значения будут сильно коррелировать, то, вероятно, модель окажется полезной для определения прочности бетона.

Напомню, что для получения корреляции между двумя числовыми векторами используется функция `cor():>`

```
cor(predicted_strength,
concrete_test$strength)      [,1][1,]  0.
8064655576
```



Не пугайтесь, если ваш результат отличается от нашего. Поскольку нейронная сеть начинает работу со случайных весов, приведенные в книге прогнозы могут быть разными для разных моделей. Если вы хотите точно сопоставить результаты, попробуйте выполнить команду `set.seed(12345)`, перед тем как начать построение нейронной сети.

Если корреляция близка к 1, это говорит о сильной линейной зависимости между двумя переменными. Следовательно, имеющая место корреляция, примерно равная 0,806, указывает на довольно сильную зависимость. Это означает, что модель достаточно хорошо работает даже с единственным скрытым узлом. Учитывая, что мы использовали только один скрытый узел, вполне вероятно, что можно улучшить эффективность модели, что мы и попробуем сделать.

Шаг 5. Повышение эффективности модели

Поскольку сети с более сложной топологией способны изучать более сложные концепции, посмотрим, что произойдет, если увеличить количество скрытых узлов до пяти. Мы, как и раньше, будем использовать функцию `neuralnet()`, но добавим параметр `hidden=5`:

```
> concrete_model2 <- neuralnet(strength ~
cement + slag + ash
+ water + superplastic
+ coarseagg +
fineagg + age, data
= concrete_train, hidden = 5)
```

Снова построив диаграмму сети (рис. 7.12), мы увидим резкое увеличение количества связей. Как это повлияло на эффективность?

```
> plot(concrete_model2)
```

Обратите внимание, что полученная ошибка (снова измеренная как SSE) уменьшилась с 5,08 в предыдущей модели до 1,63. Кроме того, количество этапов обучения возросло с 4882 до 86 849 — что неудивительно, учитывая, насколько усложнилась модель. Чем сложнее сеть, тем больше требуется итераций, чтобы найти оптимальные веса.

Применив те же шаги для сравнения прогнозируемых значений с истинными, мы получим корреляцию около 0,92, что намного лучше по сравнению с предыдущим результатом 0,80 для сети с одним скрытым узлом:

```
> model_results2 <- compute(concrete_model2,
concrete_test[1:8])> predicted_strength2 <-
model_results2$net.result>
cor(predicted_strength2,
concrete_test$strength) [ ,1][1, ]
0.9244533426
```

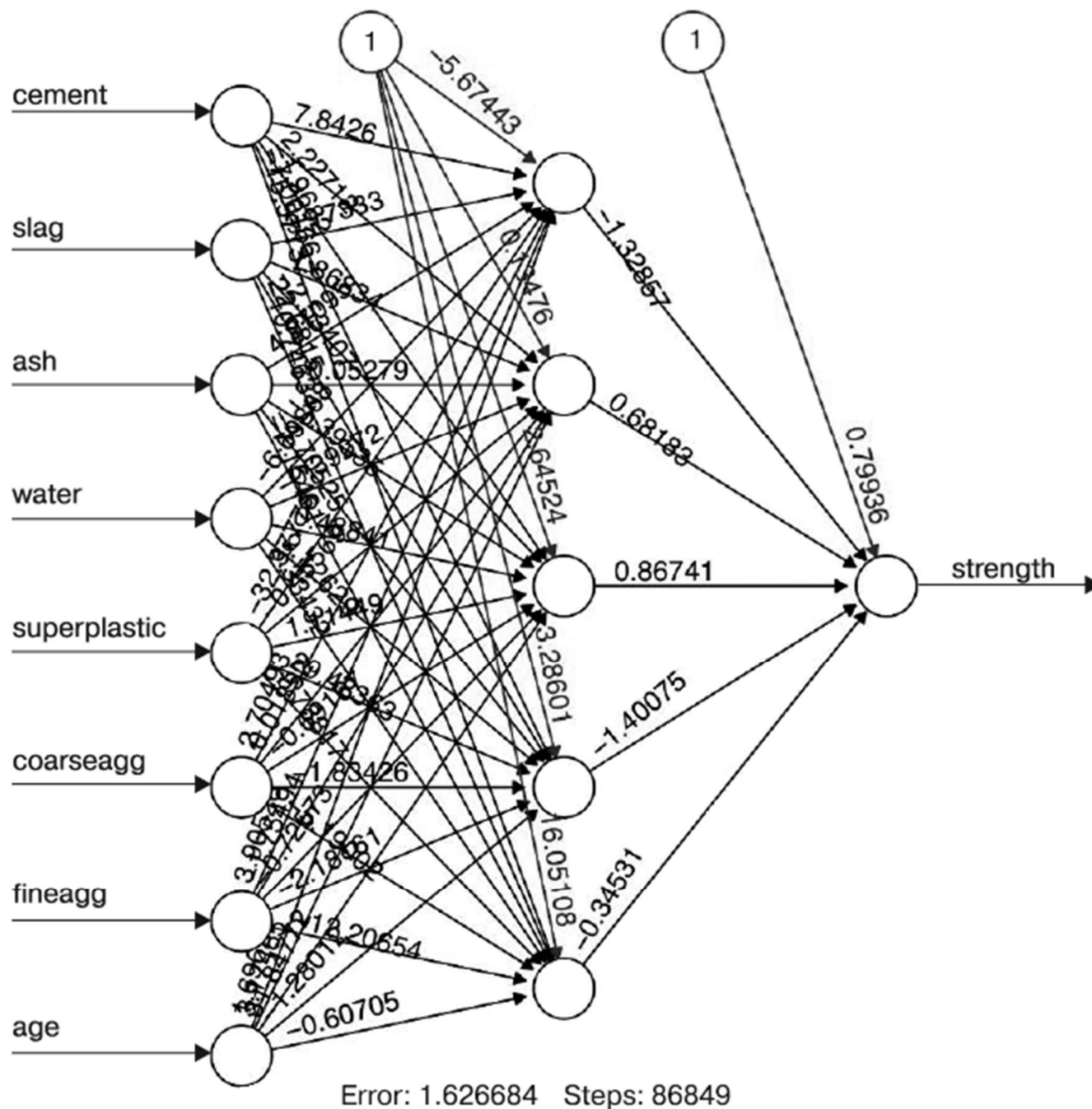


Рис. 7.12. Визуализация топологии сети с увеличенным количеством скрытых узлов

Несмотря на существенные улучшения, можно пойти еще дальше, чтобы повысить эффективность модели. В частности, есть возможность вводить дополнительные скрытые слои и изменять функцию активации сети. Внося эти изменения, мы закладываем основы для построения простейшей глубокой нейронной сети.

Выбор функции активации очень важен для глубокого обучения. Наилучшая функция для конкретной задачи обучения обычно находится экспериментальным путем, а затем широко используется сообществом исследователей машинного обучения.

В последнее время очень популярной стала функция активации, называемая *ректификационной функцией*, или *ректификатором*, — благодаря ее успешному применению в сложных задачах, таких как распознавание изображений. Узел нейронной сети, в котором в качестве функции активации используется ректификатор, называется *линейным выпрямителем* (Rectified Linear Unit, ReLU). Как показано на рис. 7.13, функция активации типа «ректификатор» описана так, что возвращает x , если x больше или равен 0, и 0 — в противном случае. Важность этой функции в том, что она, с одной стороны, нелинейная, а с другой — имеет простые математические свойства, которые делают ее недорогой в вычислительном отношении и высокоэффективной для градиентного

спуска. К сожалению, при $x = 0$ производная ректификатора не определена, поэтому ректификатор нельзя использовать в сочетании с функцией `neuralnet()`.

Вместо этого можно использовать сглаженную аппроксимацию ReLU, которая называется *softplus* или *SmoothReLU*, — функцию активации, определенную как $\log(1 + e^x)$. Как показано на рис. 7.13, функция *softplus* близка к нулю для значений x меньше 0 и приблизительно равна x для x больше 0.

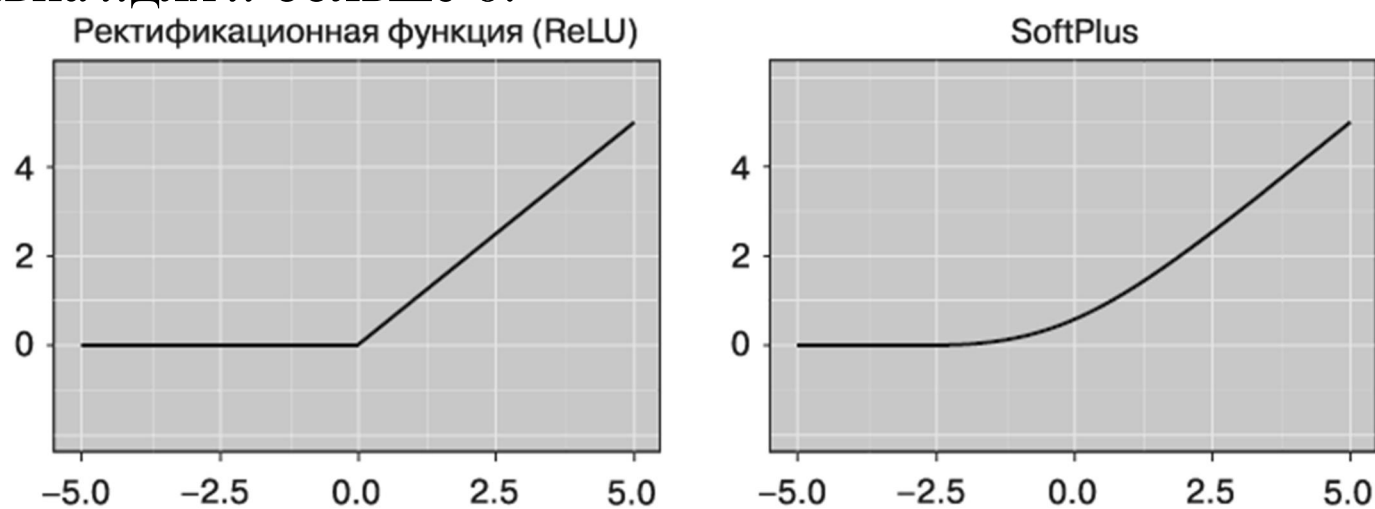


Рис. 7.13. Функция активации *softplus* — сглаженная дифференцируемая аппроксимация функции ReLU

Чтобы определить функцию `softplus()` в R, воспользуемся следующим кодом:

```
> softplus <- function(x) { log(1 + exp(x)) }
```

Такую функцию активации можно предоставить на вход `neuralnet()` с помощью параметра `act.fct`. Кроме того, добавим второй скрытый слой, состоящий из пяти узлов, присвоив параметру `hidden` значение целочисленного вектора `c(5, 5)`. В результате получим двухслойную сеть, каждый из слоев которой имеет пять узлов, и все они используют функцию активации *softplus*:

```
> set.seed(12345) > concrete_model3 <-  
neuralnet(strength ~ cement + slag  
+ ash + water +  
superplastic  
+ coarseagg +  
fineagg + age, data =  
concrete_train, hidden =  
den = c(5,  
5), act.fct =  
softplus)
```

Как и прежде, сеть можно визуализировать (рис. 7.14):

```
> plot(concrete_model3)
```

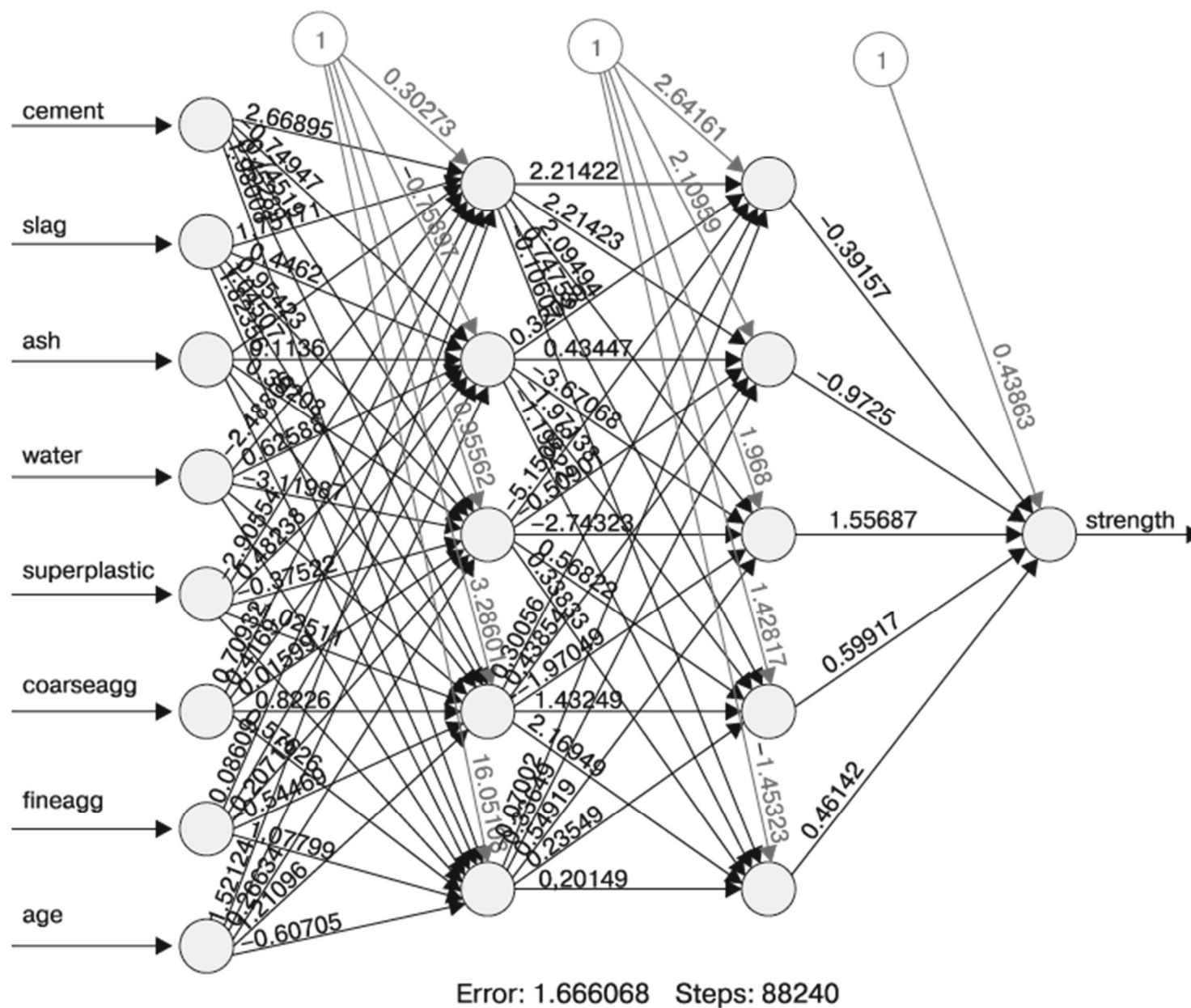


Рис. 7.14. Визуализация сети с двумя слоями скрытых узлов и с использованием функции активации softplus

Корреляцию между прогнозируемой и фактической прочностью бетона можно рассчитать следующим образом:

```
> model_results3 <- compute(concrete_model3,
concrete_test[1:8])> predicted_strength3 <-
model_results3$net.result>
cor(predicted_strength3,
concrete_test$strength)
[ ,1][1, ]
0.9348395359
```

Корреляция между прогнозируемой и фактической прочностью составила 0,935, что является лучшим из полученных до сих пор показателей. Интересно, что в оригинальной публикации Йе сообщил о корреляции 0,885. Это означает, что, приложив относительно небольшие усилия, мы смогли получить сопоставимый результат и даже превзойти результаты работы эксперта в данной области. Правда, результаты Йе были опубликованы в 1998 году, что дало нам фору более чем 20 лет дополнительных исследований в области нейронных сетей!

Следует учитывать еще одну важную деталь: поскольку мы нормализовали данные до обучения модели, прогнозы также находятся в нормализованном интервале от 0 до 1. Например, в следующем коде показан фрейм данных, построчно сравнивающий значения прочности бетона из исходного набора данных с соответствующими прогнозами:

```
> strengths <- data.frame( actual =
concrete$strength[774:1030], pred =
predicted_strength3 )> head(strengths, n =
3) actual pred774 30.14 0.286063909
1775 44.40 0.4777304648776 24.50 0.2840964250
```

Исследуя корреляцию, мы видим, что выбор нормализованных или ненормализованных данных не влияет на вычисленную статистику эффективности — точно так же, как и раньше, корреляция составляет 0,935:

```
> cor(strengths$pred, strengths$actual)[1]
0.9348395359
```

Но если бы мы вычислили другой показатель эффективности, например абсолютную разность между прогнозируемыми и фактическими значениями, то выбор масштаба был бы очень важен.

С учетом этого можно создать функцию `unnormalize()`, которая бы выполняла процедуру, обратную минимаксной нормализации, и позволяла бы преобразовывать нормализованные прогнозы в исходный масштаб:

```
> unnormalize <- function(x) { return((x *
(max(concrete$strength)) -
min(concrete$strength)) +
min(concrete$strength)) }
```

После применения написанной нами функции `unnormalize()` к прогнозам становится ясно, что масштаб новых прогнозов аналогичен исходным значениям прочности бетона. Это позволяет вычислить осмысленное значение абсолютной ошибки. Кроме того, корреляция между ненормализованными и исходными значениями прочности остается неизменной:

```
> strengths$pred_new <-
unnormalize(strengths$pred)> strengths$error <-
strengths$pred_new - strengths$actual>
head(strengths, n =
3) actual pred pred_new
error774 30.14 0.2860639091 23.62887889 -
6.511121108775 44.40 0.4777304648 39.4605363
9 -
4.939463608776 24.50 0.2840964250 23.4663647
0 -1.033635298> cor(strengths$pred_new,
strengths$actual)[1] 0.9348395359
```

Применяя нейронные сети к своим проектам, вам необходимо выполнить аналогичную последовательность шагов, чтобы вернуть данные к исходному масштабу.

Возможно, вы также обнаружите, что нейронные сети быстро становятся все более сложными, поскольку применяются для все более трудных задач обучения. Например, вы можете столкнуться с так называемой проблемой «исчезающего» малого градиента и тесно связанной с ней проблемой «взрывающегося» градиента, когда алгоритм обратного распространения ошибки не находит полезного решения, так как не сходится за разумное время. Для решения этих проблем можно попытаться изменить количество скрытых узлов, применить различные функции активации, такие как ReLU, настроить скорость обучения и т.д. На странице справки для функции `neuralnet` вы найдете дополнительную информацию о различных параметрах, которые можно настроить. Однако это приводит к другой проблеме, когда узким местом при построении высокоэффективной модели становится проверка большого количества параметров. Такова цена использования нейронных сетей и тем более сетей глубокого обучения: их огромный потенциал требует больших затрат времени и вычислительной мощности.



Подобно тому как это часто бывает в жизни, в ML можно обменять время на деньги. Использование платных ресурсов облачных вычислений, таких как Amazon Web Services (AWS) и Microsoft Azure, позволяет строить более сложные модели или быстрее тестировать многие модели. Подробнее об этом читайте в главе 12.

Метод опорных векторов

Метод опорных векторов (Support Vector Machine, SVM) можно представить как поверхность, которая образует границу между точками данных, нанесенными на график в многомерном пространстве, описывающем примеры и значения их признаков. Цель SVM состоит в том, чтобы построить плоскую границу — *гиперплоскость*, которая бы делила пространство таким образом, чтобы по обеим ее сторонам образовались однородные группы. Таким образом, в обучении SVM сочетаются аспекты как обучения методом ближайших соседей на основе экземпляров, описанного в главе 3, так и моделирования методом линейной регрессии, рассмотренного в главе 6. Это чрезвычайно мощное сочетание, позволяющее SVM моделировать очень сложные взаимосвязи.

Несмотря на то что базовая математика, лежащая в основе SVM, существовала десятки лет, интерес к этим методам значительно вырос после того, как они стали применяться для ML. Популярность этих

методов возросла после громких историй успеха в решении сложных задач обучения, а также после разработки алгоритмов SVM, которые были отмечены наградами и реализованы в хорошо поддерживаемых библиотеках на многих языках программирования, включая R. После этого методы SVM были приняты широкой аудиторией. В противном случае, вероятно, невозможно было бы применять сложную математику, необходимую для реализации SVM. Положительным моментом является то, что хотя математика, возможно, и сложная, основные концепции вполне понятны.

Методы SVM могут быть адаптированы для использования практически любых типов задач обучения, включая классификацию и числовое прогнозирование. Многие ключевые успехи этого алгоритма относятся к распознаванию образов. В число наиболее известных областей применения этих методов входят следующие:

- классификация данных по экспрессии генов микрочипов в биоинформатике для выявления рака и других генетических заболеваний;
- текстовая категоризация, такая как определение языка, используемого в документе, или классификация документов по теме;
- обнаружение редких, но важных событий, таких как выход из строя двигателя внутреннего сгорания, нарушение безопасности или землетрясение.

Методы SVM легче всего понять на примере бинарной классификации — именно так они обычно применяются. Поэтому в остальных разделах мы сосредоточимся только на SVM-классификаторах. Принципы, подобные представленным здесь, используются также при адаптации SVM-методов для числового прогнозирования.

Классификация гиперплоскостями

Как отмечалось ранее, в SVM-методах строится граница, называемая гиперплоскостью, для разделения данных на группы с одинаковыми значениями классов. Например, на рис. 7.15 показаны гиперплоскости, которые разделяют группы кругов и квадратов в двух и трех измерениях. Поскольку группы кругов и квадратов можно идеально разделить прямой линией или плоскостью, то такие группы называются *линейно разделимыми*. Сначала рассмотрим простой случай, когда это условие выполняется. Однако SVM также можно распространить на те задачи, в которых точки не являются линейно разделимыми.

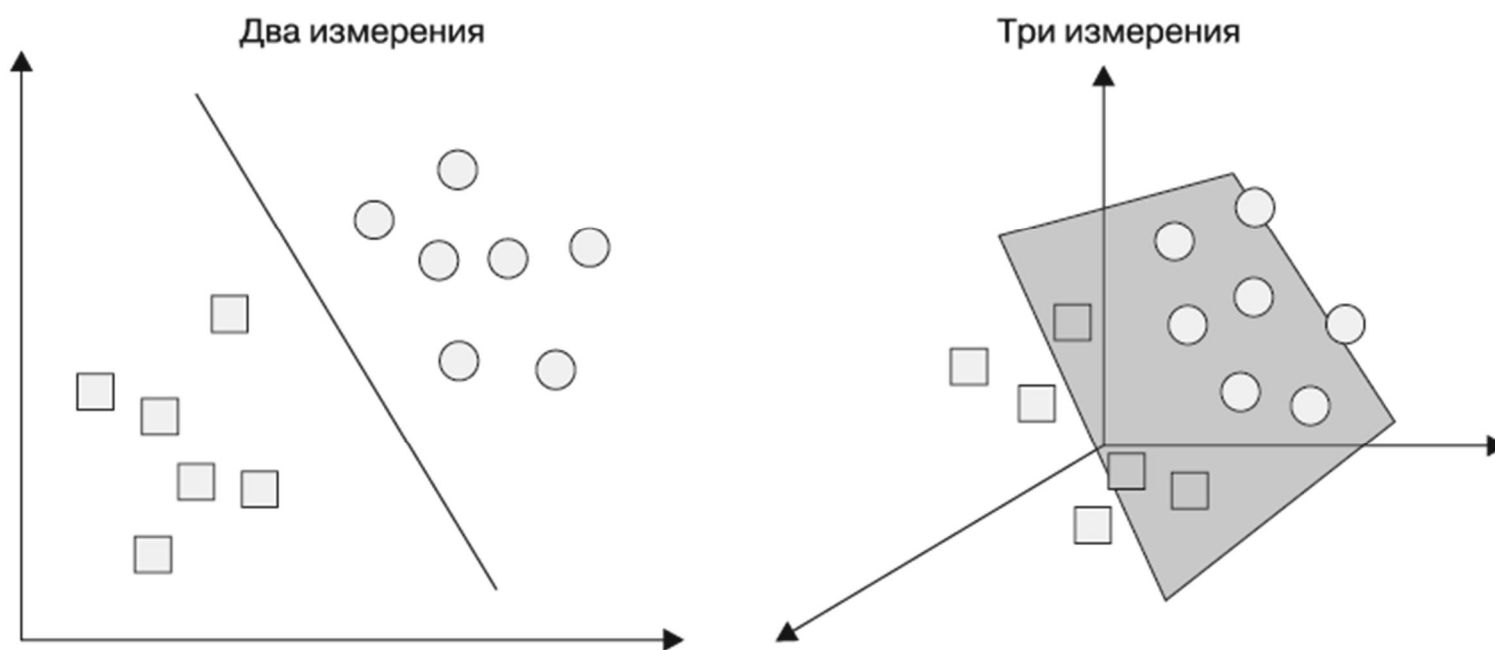


Рис. 7.15. Квадраты и круги линейно разделимы в двух и трех измерениях



Для удобства гиперплоскость обычно изображается в виде прямой в 2D-пространстве, но это делается потому, что трудно проиллюстрировать пространство, состоящее из более чем двух измерений. На самом деле гиперплоскость — это плоскость многомерного пространства, концепция, восприятие которой может оказаться трудным для понимания.

В двух измерениях задача алгоритма SVM состоит в том, чтобы найти линию, разделяющую два класса. Как показано на рис. 7.16, существует несколько возможных разделительных линий между группами кругов и квадратов. Три такие возможности обозначены на рис. 7.16 буквами *a*, *b* и *c*. Какую из них выберет алгоритм?

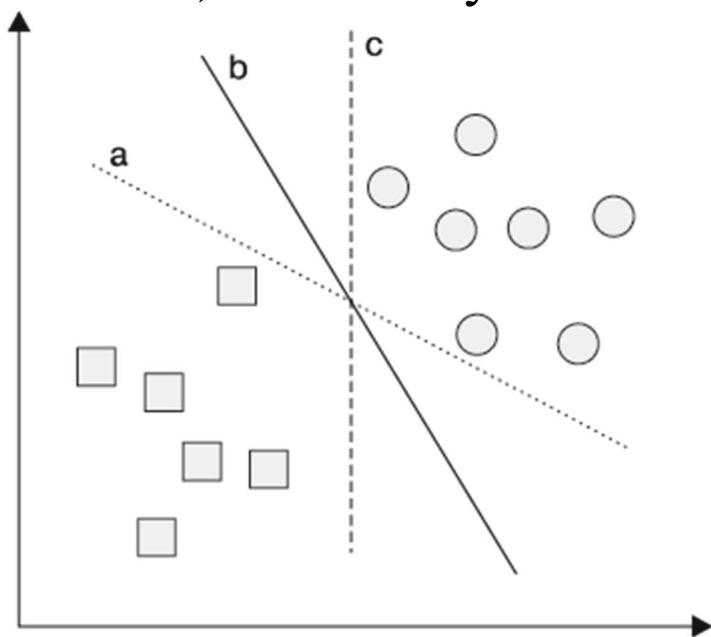


Рис. 7.16. Три линии из множества возможных, разделяющие квадраты и круги

Ответ на этот вопрос включает в себя поиск *оптимально разделяющей гиперплоскости* (Maximum Margin Hyperplane, ММН), которая лучше всего разделяет два класса. Любая из трех линий, разделяющих окружности и квадраты, правильно классифицирует все точки данных, однако линия, которая приводит к наибольшему разделению, является наилучшим обобщением для будущих данных. Максимальный запас повышает вероятность того, что, даже если появится случайный шум, каждый класс останется на своей стороне границы.

Опорные векторы (на рис. 7.17 обозначены стрелками) являются точками из каждого класса, ближайшими к ММН. Каждый класс должен иметь хотя бы один опорный вектор, но их может быть и несколько. Только опорные векторы определяют ММН. Это главная особенность SVM; благодаря опорным векторам способ хранения модели классификации является очень компактным, даже если признаков чрезвычайно много.

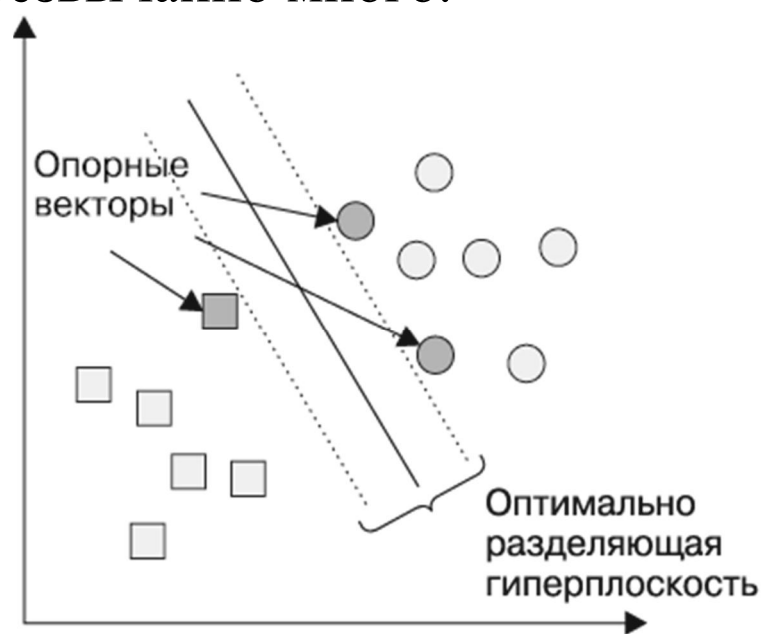


Рис. 7.17. Оптимально разделяющая гиперплоскость определяется опорными векторами

Алгоритм определения опорных векторов основан на векторной геометрии и состоит из довольно сложных математических операций, которые в данной книге не рассматриваются. Однако основные принципы этого процесса довольно просты.



Подробнее о математических основах SVM читайте в классической статье: Cortes C., Vapnik V. Support-Vector Networks. Machine Learning, 1995. Vol. 20. P. 273–297. Обсуждение этого вопроса на уровне для начинающих можно найти в статье: Bennett K.P., Campbell C. Support Vector Machines: Hype or Hallelujah? SIGKDD Explorations, 2000. Vol. 2.P. 1–13, а более подробное описание — в книге: Steinwart I., Christmann A. Support Vector Machines. New York: Springer, 2008.

Случай линейно разделимых данных

Чтобы найти оптимально разделяющую плоскость, проще всего исходить из предположения, что классы являются линейно разделимыми. В этом случае ММН располагается на максимальном расстоянии от внешних границ двух групп точек данных. Эти внешние границы называются *выпуклой оболочкой*. Тогда ММН представляет собой серединный перпендикуляр самого короткого отрезка, соединяющего две выпуклые оболочки (рис. 7.18). Сложные компьютерные алгоритмы, использующие метод *квадратичной оптимизации*, находят оптимально разделяющую гиперплоскость именно таким способом.

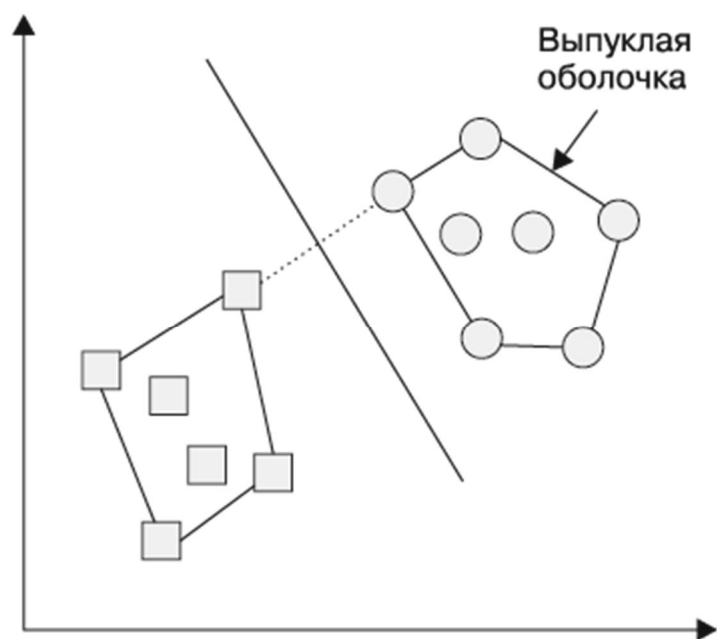


Рис. 7.18. ММН — это серединный перпендикуляр кратчайшего пути между выпуклыми оболочками

Еще один (эквивалентный) подход включает в себя поиск в пространстве всех возможных гиперплоскостей двух параллельных плоскостей, которые делят точки на однородные группы, но сами расположены как можно дальше друг от друга. Можно использовать такую метафору: этот процесс похож на попытку найти самый толстый матрас, который можно поднять по лестнице в вашу спальню.

Чтобы понять, как происходит процесс поиска, нам нужно точно определить, что мы подразумеваем под гиперплоскостью. В n -мерном пространстве используется следующее уравнение:

$$\vec{w} \cdot \vec{x} + b = 0.$$

Возможно, вам не знакомы эти обозначения: стрелки над буквами указывают на то, что они являются векторами, а не одиночными числами. В частности, w — это вектор из n весов, то есть $\{w_1, w_2, \dots, w_n\}$, а b — число, называемое смещением. Концептуально смещение — это то же, что сдвиг в уравнении с угловым коэффициентом, которое обсуждалось в главе 6.



Если вам трудно представить плоскость в многомерном пространстве, не беспокойтесь о деталях. Просто представьте уравнение как способ описания поверхности, очень похожий на уравнение с угловым коэффициентом ($y = mx + b$), которое используется для описания прямых в двумерном пространстве.

Цель процесса — с помощью этой формулы найти набор весов, которые бы определяли две гиперплоскости, следующим образом:

$$\vec{w} \cdot \vec{x} + b \geq +1;$$

$$\vec{w} \cdot \vec{x} + b \leq -1.$$

Необходимо, чтобы эти гиперплоскости были определены так, чтобы все точки одного класса располагались над первой гиперплоскостью, а все точки второго класса — под второй гиперплоскостью. Это возможно при условии, что данные являются линейно разделимыми.

В векторной геометрии расстояние между этими двумя плоскостями определяется следующим образом:

$$\frac{2}{\|\bar{w}\|}.$$

Здесь $\|w\|$ — евклидова норма (расстояние от начала координат до вектора w). Поскольку $\|w\|$ является знаменателем, то, чтобы получить максимальное расстояние, нужно минимизировать $\|w\|$. Поэтому задача обычно переформулируется как набор ограничений следующим образом:

$$\min \frac{1}{2} \|\bar{w}\|^2$$

если $(\bar{w} \cdot \bar{x}_i - b) \geq 1, \forall \bar{x}_i.$

Это кажется запутанным, однако на самом деле это не так уж сложно понять. По сути, первая строка подразумевает, что нам нужно минимизировать евклидову норму (чтобы расчет был проще, возвести в квадрат и разделить на 2). Во второй строке отмечается, что это должно выполняться при условии (если), что каждая точка данных x_i классифицирована правильно. Обратите внимание: y_i указывает на значение класса (преобразованное в +1 или в -1), а перевернутая буква «A» является сокращением «для каждого».

Как и для других методов определения оптимально разделяющей гиперплоскости, поиск решения этой проблемы является задачей, которую лучше всего предоставить программному обеспечению квадратичной оптимизации. Несмотря на интенсивную загрузку процессора, специализированные алгоритмы способны быстро решать такие задачи даже для довольно больших наборов данных.

Случай линейно неразделимых данных

Поскольку мы уже проработали теорию, лежащую в основе SVM, может возникнуть вопрос: что происходит, если данные не являются линейно разделимыми? Решением этой проблемы является использование *ослабляющей переменной*, которая создает мягкую границу, позволяющую некоторым точкам попадать не на свою сторону. На рис. 7.19 показаны две точки, попадающие на «чужую» сторону, и соответствующие ослабляющие переменные (обозначенные греческой буквой «кси»):

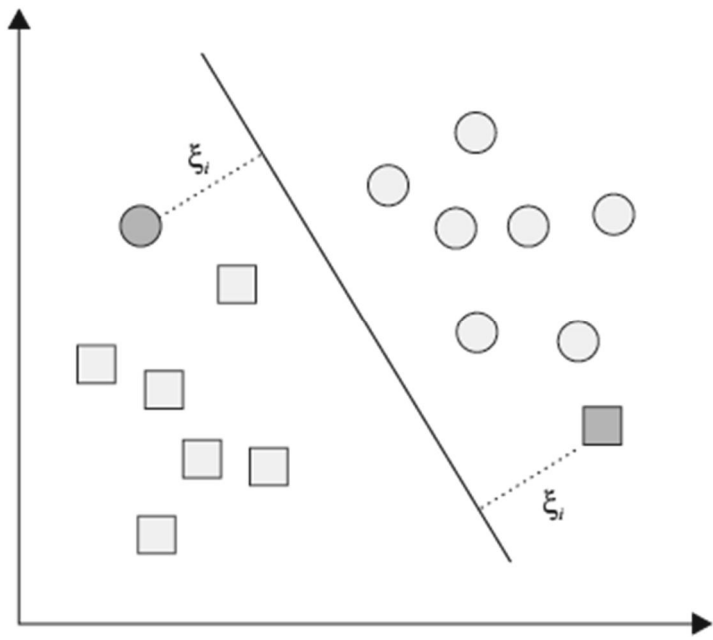


Рис. 7.19. К точкам, попадающим не на свою сторону от границы, применяется штраф

Штраф (cost value) (обозначается буквой C) применяется ко всем точкам, которые нарушают ограничения, и вместо того, чтобы искать оптимально разделяющую гиперплоскость, алгоритм старается минимизировать общий штраф. Поэтому можно переформулировать проблему оптимизации следующим образом:

$$\min \frac{1}{2} \|\bar{w}\|^2 + C \sum_{i=1}^n \xi_i$$

если $y_i (\bar{w} \cdot \bar{x} - b) \geq 1 - \xi_i, \forall \bar{x}_i, \xi_i \geq 0$.

Если вы запутались, не волнуйтесь — вы не одиноки. К счастью, SVM-пакеты оптимизируют все это, так что от вас не требуется разбираться в технических деталях. Важно лишь понимать, что означает добавление параметра штрафа C . Изменение этого значения приводит к корректировке штрафа для точек, не попадающих на свою сторону гиперплоскости. Чем больше штраф, тем сложнее будет алгоритму оптимизации достичь 100%-ного разделения. Однако более низкие значения затрат требуют более широкой общей разделяющей границы. Важно найти баланс между этими двумя параметрами, чтобы получить модель, которая будет хорошо обобщать будущие данные.

Использование ядер в нелинейных пространствах

Во многих реальных наборах данных зависимости между переменными являются нелинейными. Как вы только что убедились, метод SVM можно обучать и на таких данных путем добавления ослабляющей переменной, что допускает ошибочную классификацию некоторых примеров. Однако это не единственный способ решить проблему нелинейности. Главная особенность SVM-методов — их способность отображать задачу в пространстве с увеличенным числом измерений, используя процесс, известный как *трюк с ядром* (kernel trick). При этом нелинейные отношения могут внезапно оказаться линейными.

Возможно, это кажется бессмысленным, однако на самом деле это довольно легко проиллюстрировать на примере. На рис. 7.20 диаграмма

рассеяния слева отражает нелинейную связь между классом погоды — ясно или снежно — и двумя объектами — широтой и долготой. Точки в центре диаграммы являются членами класса «снежно», а точки на краях относятся к классу «ясно». Такие данные могли быть получены из метеорологических отчетов, поступающих с метеостанций, расположенных у вершины горы, и со станций, находящихся у подножия.

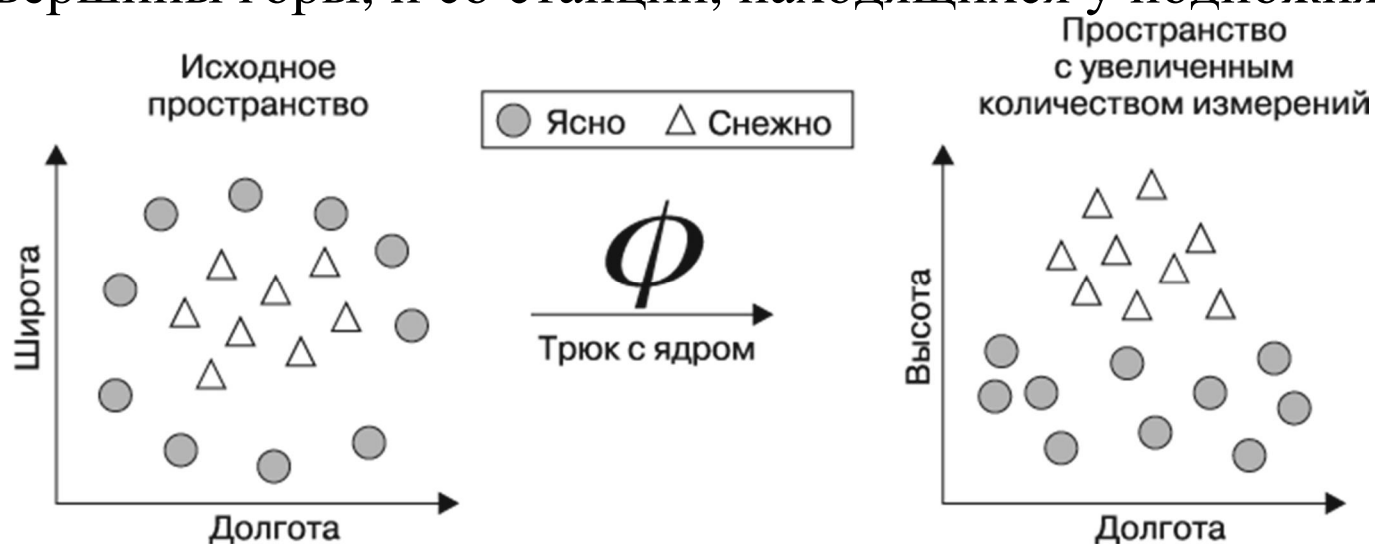


Рис. 7.20. Трюк с ядром может помочь преобразовать нелинейную задачу в линейную

В правой части рисунка изображена диаграмма после применения трюка с ядром. Здесь мы смотрим на данные через призму нового измерения — высоты. С добавлением этого признака классы становятся идеально линейно разделимыми. Это стало возможным, потому что мы по-новому рассмотрели данные. На диаграмме слева мы видим гору с высоты птичьего полета, а справа — ту же гору на уровне земли, но с некоторого расстояния. Теперь тенденция очевидна: снежная погода чаще бывает на больших высотах.

Таким образом, чтобы создать разделение, SVM-методы с нелинейными ядрами добавляют дополнительные измерения к данным. По сути, трюк с ядром означает создание новых признаков, которые отражают математические отношения между измеряемыми характеристиками.

В частности, признак высоты над уровнем моря может быть математически описан как зависимость между широтой и долготой: чем ближе точка к центру каждой из этих осей, тем больше высота. Это позволяет SVM-методам изучать понятия, которые не содержались в исходных данных в виде явно измеренных значений.

SVM-методы с нелинейными ядрами являются очень мощными классификаторами, хотя и имеют ряд недостатков, как показано в табл. 7.2.

Таблица 7.2

Преимущества	Недостатки
<p>Метод может использоваться для задач классификации и числового прогнозирования.</p> <p>Не слишком подвержен влиянию зашумленных данных и переобучению.</p> <p>Иногда проще в использовании, чем нейронные сети, особенно благодаря существованию нескольких хорошо поддерживаемых SVM-алгоритмов.</p> <p>Получил популярность благодаря высокой точности и резонансным победам в соревнованиях по интеллектуальному анализу данных</p>	<p>Поиск наилучшей модели требует проверки различных сочетаний ядер и параметров модели.</p> <p>Может оказаться медленным, особенно если набор входных данных имеет большое количество признаков или примеров.</p> <p>В результате получается сложная модель типа «черный ящик», которую трудно, если вообще возможно, интерпретировать</p>

Как правило, функции ядра имеют следующую форму. Функция, обозначаемая греческой буквой «фи» — $\phi(x)$, — представляет собой отображение данных в другое пространство. Таким образом, обобщенная функция ядра применяет некоторое преобразование к векторам признаков x_i и x_j и объединяет их с помощью *скалярного произведения*, которое принимает два вектора и возвращает число.

$$K(\bar{x}_i, \bar{x}_j) = \phi(\bar{x}_i) \cdot \phi(\bar{x}_j).$$

С помощью этой формы были разработаны функции ядра для многих областей применения. Некоторые из наиболее часто используемых функций ядра перечислены ниже. Почти все пакеты программного обеспечения для SVM включают в себя эти ядра, а также многие другие.

Линейное ядро не преобразует данные вообще. Таким образом, его можно описать просто как скалярное произведение признаков:

$$K(\bar{x}_i, \bar{x}_j) = \bar{x}_i \cdot \bar{x}_j.$$

Полиномиальное ядро степени d добавляет простое нелинейное преобразование данных:

$$K(\bar{x}_i, \bar{x}_j) = (\bar{x}_i \cdot \bar{x}_j + 1)^d.$$

Сигмоидное ядро приводит к SVM-модели, в чем-то подобной нейронной сети с сигмоидной функцией активации. Греческие буквы «каппа» и «дельта» обозначают параметры ядра:

$$K(\bar{x}_i, \bar{x}_j) = \tanh(k\bar{x}_i \cdot \bar{x}_j - \delta).$$

Гауссово RBF-ядро похоже на нейронную сеть RBF. RBF-ядро хорошо работает для многих типов данных и считается отправной точкой для многих задач обучения:

$$K(\bar{x}_i, \bar{x}_j) = e^{-\frac{\|\bar{x}_i - \bar{x}_j\|^2}{2\sigma^2}}.$$

Не существует надежного правила выбора ядра для конкретной задачи обучения. Подбор ядра в значительной степени зависит от концепции, которую необходимо изучить, а также от количества тренировочных данных и взаимосвязи между признаками. Зачастую обучение и оценка на тестовых данных нескольких SVM-методов происходят путем проб и ошибок. При этом во многих случаях выбор ядра является произвольным, поскольку эффективность может отличаться лишь незначительно. Чтобы увидеть, как это работает на практике, применим наши знания о SVM-классификации к реальной задаче.

Пример: оптическое распознавание символов с помощью SVM

Для многих типов алгоритмов машинного обучения обработка изображений является сложной задачей. Взаимосвязи между паттернами пикселей и более высокими концепциями очень сложны, и их трудно распознавать. Например, человек легко различает лицо, кошку или букву «А», но описать эти закономерности в строгих правилах сложно. Кроме того, исходные изображения часто бывают зашумленными. Одно и то же изображение может незначительно различаться в зависимости от того, как оно было получено: при разном освещении, ориентации в пространстве и расположении объекта.

SVM-методы хорошо подходят для решения проблем, связанных с данными, представленными в виде изображений. Эти методы, способные к обучению сложным паттернам без чрезмерной чувствительности к шуму, позволяют обнаруживать визуальные закономерности с высокой степенью точности. Более того, главный недостаток SVM — представление модели в виде «черного ящика» — при обработке изображений не так уж критичен. Если SVM-метод позволяет отличить кошку от собаки, то не имеет значения, как именно он это делает.

В этом разделе построим модель, аналогичную той, что используется в ядре программного обеспечения для *оптического распознавания символов* (Optical Character Recognition, OCR), часто поставляемого в комплекте с настольными сканерами или используемого в приложениях для смартфонов. Целью такого программного обеспечения является обработка бумажных документов путем преобразования печатного или рукописного текста в электронную форму, пригодную для последующего сохранения в базе данных.

Конечно, из-за множества вариантов рукописного ввода и печатных шрифтов это весьма сложная задача. Однако пользователи программного обеспечения ожидают совершенства, поскольку в бизнесе ошибки или опечатки чреватые неприятными последствиями. Проверим, подходит ли наш SVM-метод для решения этой задачи.

Шаг 1. Сбор данных

Когда *OCR-программа* впервые обрабатывает документ, она делит лист на прямоугольники, где каждый содержит один *глиф* — этим термином обозначают букву, символ или цифру. Затем программа попытается сопоставить глиф, находящийся в каждой ячейке, со множеством всех символов, которые она способна распознать. Затем отдельные символы могут быть объединены в слова, которые при желании можно проверить по словарю на языке документа.

Для этого упражнения предположим, что мы уже разработали алгоритм для разделения документа на прямоугольные области, каждая из которых содержит один глиф. Мы также предположим, что документ содержит только буквы английского языка. Поэтому построим модель процесса, который включает в себя сопоставление глифов с одной из 26 букв от A до Z.

Для этого воспользуемся набором данных, переданным У. Фреем (W. Frey) и Д. Дж. Слейтом (D. J. Slate) в репозиторий машинного обучения UCI Machine Learning Repository (<http://archive.ics.uci.edu/ml>). Этот набор данных состоит из 20 000 примеров, содержащих 26 заглавных букв английского алфавита и напечатанных 20 различными вариантами с использованием произвольно измененных и искаженных черно-белых шрифтов.



Подробнее об этих данных читайте в публикации: Slate D.J., Frey P.W. Letter Recognition Using Holland-Style Adaptive Classifiers. Machine Learning, 1991. Vol. 6.P. 161–182.

На рис. 7.21, опубликованном Фреем и Слейтом, показаны некоторые напечатанные глифы. Искаженные таким образом буквы вызывают трудности при компьютерной идентификации, но легко распознаются человеком.

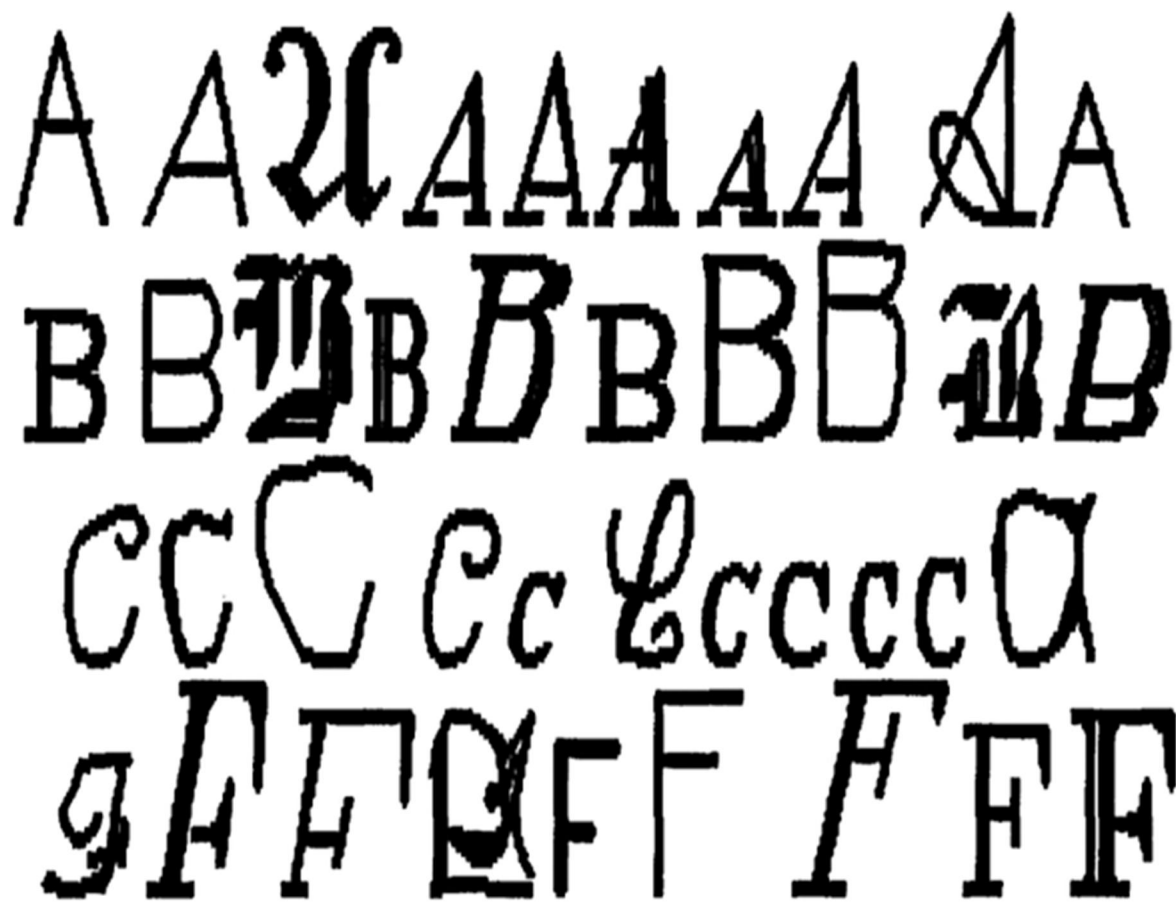


Рис. 7.21. Примеры глифов, которые SVM-алгоритм попытается распознать

Шаг 2. Исследование и подготовка данных

Согласно документации, предоставленной Фреем и Слейтом, при сканировании глифов в компьютер они преобразуются в пикселы и для каждого глифа записываются 16 статистических атрибутов.

Эти атрибуты измеряют такие характеристики, как горизонтальные и вертикальные размеры глифа; соотношение черных (по сравнению с белыми) пикселей; среднее горизонтальное и вертикальное положение пикселей. Предполагается, что различия в концентрации черных и белых пикселей на разных участках поля должны позволить различать 26 букв алфавита.



Для того чтобы выполнить этот пример, загрузите файл `letterdata.csv` и сохраните его в рабочем каталоге `R`.

Когда данные будут считаны в среду `R`, нам надо убедиться, что мы получили данные, в которых каждый пример содержит 16 признаков, определяющих класс буквы `letter`. Как и ожидалось, этот класс имеет 26 уровней:

```
> letters <- read.csv("letterdata.csv")>
str(letters)'data.frame':    20000 obs. of 17
variables:$ letter: Factor w/ 26 levels
"A","B","C","D",...$ xbox  : int  2 5 4 7 2 4 4 1
2 11 ...$ ybox  : int  8 12 11 11 1 11 2 1 2 15
...$ width  : int  3 3 6 6 3 5 5 3 4 13 ...$
height: int  5 7 8 6 1 8 4 2 4 9 ...$ onpix :
int  1 2 6 3 1 3 4 1 2 7 ...$ xbar  : int  8 10
10 5 8 8 8 8 10 13 ...$ ybar  : int  13 5 6 9 6 8
7 2 6 2 ...$ x2bar  : int  0 5 2 4 6 6 6 2 2 6
...$ y2bar  : int  6 4 6 6 6 9 6 2 6 2 ...$ xybar
: int  6 13 10 4 6 5 7 8 12 12 ...$ x2ybar:
int  10 3 3 4 5 6 6 2 4 1 ...$ xy2bar: int  8 9 7
10 9 6 6 8 8 9 ...$ xedge  : int  0 2 3 6 1 0 2 1
1 8 ...$ xedgey: int  8 8 7 10 7 8 8 6 6 1 ...$
yedge  : int  0 4 3 2 5 9 7 2 1 1 ...$ yedgex:
int  8 10 9 8 10 7 10 7 7 8 ...
```

Обучаемый SVM-алгоритм требует, чтобы все признаки были числовыми и, кроме того, чтобы каждый признак масштабировался к довольно небольшому интервалу. В данном случае все признаки являются целыми числами, поэтому нам не нужно преобразовывать какие-либо из них в числа. Однако для некоторых признаков диапазоны этих целочисленных переменных довольно широки. Это значит, что нужно нормализовать или стандартизировать эти данные. Однако пока можно пропустить этот шаг, поскольку `R`-пакет, который будет использоваться для построения SVM-модели, автоматически изменит масштаб.

Учитывая, что подготовка данных не требуется, можно сразу перейти к следующим этапам машинного обучения — собственно обучению и тестированию. В предыдущих задачах анализа мы случайным образом делили данные на тренировочный и тестовый наборы. Мы могли бы так поступить и здесь, но Фрей и Слейт уже рандомизировали данные и предлагают использовать первые 16 000 записей (80 %) для построения модели и остальные 4000 записей (20 %) — для тестирования. Следуя их советам, мы можем создать тренировочные и тестовые фреймы данных следующим образом:

```
> letters_train <- letters[1:16000, ]>
letters_test  <- letters[16001:20000, ]
```

Итак, данные готовы к работе. Приступим к построению классификатора.

Шаг 3. Обучение модели на данных

В R есть несколько выдающихся пакетов для построения SVM-моделей. Пакет `e1071`, разработанный кафедрой статистики Венского технического университета (TU Wien), предоставляет R-интерфейс для отмеченной наградами библиотеки LIBSVM — широко используемой SVM-программы с открытым исходным кодом, написанной на C++. Если вы уже знакомы с LIBSVM, то можете начать с нее.



Подробнее об LIBSVM читайте на сайте авторов этой библиотеки по адресу <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.

Аналогично, если вы уже освоили алгоритм SVMlight, можете воспользоваться пакетом `klaR`, разработанным на кафедре статистики Дортмундского технологического университета (TU Dortmund). Этот пакет предоставляет функции для работы с данной реализацией SVM непосредственно из R.



Подробнее об SVMlight читайте на сайте <http://svmlight.joachims.org/>.

Наконец, если вы начинаете с нуля, то лучше всего воспользоваться SVM-функциями из пакета `kernlab`. Преимуществом этого пакета является то, что он изначально разрабатывался не на C или C++, а на R, что позволяет легко его настраивать; у него нет каких-либо внутренних, скрытых элементов. Пожалуй, еще более важно, что `kernlab`, в отличие от других пакетов, можно использовать совместно с пакетом `caret`, который позволяет обучать и оценивать SVM-модели с помощью различных автоматизированных методов (см. главу 11).



Подробное введение в kernlab вы найдете в статье по адресу <http://www.jstatsoft.org/v11/i09/>.

Синтаксис для обучения SVM-классификаторов с помощью kernlab следующий. Если вы будете использовать другой пакет, то его команды во многом похожи. По умолчанию функция ksvm() использует гауссово RBF-ядро, но предоставляются и некоторые другие варианты.

Синтаксис метода опорных векторов

Использование функции ksvm() из пакета kernlab

Построение модели:

```
m <- ksvm(target ~ predictors, data = mydata,  
          kernel = "rbfdot", c = 1)
```

target – модель, которая будет построена в результате обучения на фрейме данных mydata;

predictors – R-формула, определяющая признаки из фрейма данных mydata, которые будут использоваться при прогнозе;

data – фрейм данных, которому принадлежат target и predictors;

kernel – нелинейное отображение, такое как "rbfdot" (радиально-базисная функция), "polydot" (полином), "tanhdot" (тангенса), "vanilladot" (линейная функция);

c – число, определяющее штраф при нарушении ограничений: насколько высокую цену придется заплатить за «мягкое» значение, тем уже граница.

Функция возвращает SVM-объект, который можно использовать для прогнозирования.

Прогнозирование:

```
p <- ksvm(m, test, type = "response")
```

m – модель, обученная с помощью функции ksvm();

test – фрейм данных, содержащий тестовые данные с теми же признаками, что и у тренировочных данных, использованных для обучения;

type определяет, должны ли прогнозы предсказывать класс точно ("response") или с определенной вероятностью ("probabilities") для каждого уровня класса.

Функция возвращает вектор (или матрицу) прогнозируемых классов (или их вероятностей), в зависимости от значения параметра type.

Пример:

```
letter_classifier <- ksvm(letter ~ ., data = letters_train,  
                        kernel = "vanilladot")  
letter_prediction <- predict(letter_classifier, letters_test)
```

Для того чтобы обеспечить основу для оценки эффективности SVM-метода, начнем с обучения простого линейного SVM-классификатора. Если вы этого еще не сделали, установите пакет kernlab в свою библиотеку с помощью команды `install.packages("kernlab")`. После этого можно вызвать функцию `ksvm()` для тренировочных данных и выбрать линейное (vanilla) ядро, используя опцию `vanilladot`:

```
> library(kernlab)> letter_classifier <-  
ksvm(letter ~ ., data =
```

```
letters_train, kernel
= "vanilladot")
```

Эта операция может занять какое-то время, зависящее от мощности компьютера. После завершения операции введите имя сохраненной модели, чтобы увидеть основную информацию о параметрах обучения и соответствии модели:

```
> letter_classifierSupport Vector Machine
object of class "ksvm"SV type: C-
svc      (classification)  parameter : cost C =
1Linear (vanilla) kernel function.Number of
Support Vectors : 7037Objective Function Value :
-14.1746 -20.0072 -23.5628 -6.2009 -7.5524-
32.7694 -49.9786 -18.1824 -62.1111 -32.7284 -
16.2209...Training error : 0.130062
```

Эта информация очень мало говорит о том, насколько модель действительно эффективна. Для того чтобы узнать, насколько хорошо она обобщает новые данные, нужно проверить ее на тестовом наборе данных.

Шаг 4. Оценка эффективности модели

Функция `predict()` позволяет использовать модель классификации букв для прогнозирования на тестовом наборе данных:

```
> letter_predictions <-
predict(letter_classifier, letters_test)
```

Поскольку мы не указали параметр `type`, по умолчанию используется `type="response"`. В этом случае функция возвращает вектор, содержащий спрогнозированную букву для каждой строки значений тестовых данных. Воспользовавшись функцией `head()`, мы видим, что первые шесть спрогнозированных букв — это U, N, V, X, N и H:

```
> head(letter_predictions)[1] U N V X N
НLevels: A B C D E F G H I J K L M N O P Q R S T
U V W X Y Z
```

Чтобы проверить, насколько хорошо работает наш классификатор, нужно сравнить прогнозируемую букву с истинной, содержащейся в наборе тестовых данных. Для этого воспользуемся функцией `table()` (здесь показана только часть таблицы):

```
> table(letter_predictions,
letters_test$letter)letter_predictions      A      B
      C      D      E      A  144      0      0
      0      0      B      0  121      0      5
      2      C      0      0  120      0      4
```

	D	2	2	0	156	0
E	0	0	5	0	127	

Значения 144, 121, 120, 156 и 127, расположенные на диагонали, указывают на общее число записей, в которых прогнозируемая буква соответствует истинному значению. Указано и количество ошибок. Например, значение 5 в строке В и столбце D говорит о том, что в пяти случаях буква D была ошибочно идентифицирована как В.

Рассматривая каждый тип ошибок по отдельности, можно выявить интересную закономерность в отношении определенных типов букв, с которыми у модели возникают проблемы, но это занимает много времени. Можно упростить оценку, рассчитав общую точность. Она учитывает только, был ли прогноз правильным или нет, и игнорирует тип ошибки.

Следующая команда возвращает вектор значений TRUE или FALSE, в зависимости от того, соответствует ли прогнозируемая буква модели фактической букве из набора тестовых данных:

```
> agreement <- letter_predictions ==
letters_test$letter
```

С помощью функции table() мы увидим, что классификатор правильно распознал букву в 3357 из 4000 тестовых записей:

```
>
table(agreement) agreement FALSE TRUE 643 3357
```

В процентном отношении точность составляет около 84 %:

```
>
prop.table(table(agreement)) agreement FALSE
TRUE 0.16075 0.83925
```

Обратите внимание, что в 1991 году, когда Фрей и Слейт опубликовали набор данных, они сообщили, что точность распознавания составила примерно 80 %. Используя всего несколько строк R-кода, мы смогли превзойти их результат — правда, у нас есть преимущество в несколько десятков лет дополнительных исследований в области машинного обучения. С учетом этого вполне вероятно, что модель можно сделать еще лучше.

Шаг 5. Повышение эффективности модели

Теперь уделим немного времени согласованию с контекстом эффективности SVM-модели, которую мы обучили для идентификации букв алфавита по данным, полученным из изображения. Модель, построенная с помощью всего одной строки R-кода, позволила достичь почти 84%-ной точности, что немного превысило контрольную точность в процентах, опубликованную академическими исследователями в 1991 году. Несмотря на то что точность в 84 % недостаточно высока, чтобы

быть полезной для реальной OCR-программы, тот факт, что относительно простая модель способна достичь этого уровня, сам по себе является замечательным достижением. Учтите, что вероятность случайного совпадения прогноза модели с действительным значением довольно мала — менее 4 %. Это означает, что наша модель работает в 20 раз лучше, чем случайное совпадение. Не менее примечательно то, что, изменяя параметры SVM-функции и немного усложнив модель обучения, можно получить модель, которая окажется полезной в реальных приложениях.



Для того чтобы рассчитать вероятность случайного совпадения прогнозов SVM-модели с фактическими значениями, применим правило совместной вероятности для независимых событий, описанное в главе 4. Поскольку в тестовом наборе данных содержится 26 букв, каждая из которых появляется примерно с одинаковой частотой, вероятность того, что любая буква будет правильно спрогнозирована, составляет $(1 / 26) \cdot (1 / 26)$. Поскольку всего букв 26, то общая вероятность совпадения составляет $26 \cdot (1 / 26) \cdot (1 / 26) = 0,0384$, то есть 3,84 %.

Изменение функции ядра в SVM-методе

В предыдущей SVM-модели использовалась простая линейная функция ядра. Применив более сложную функцию ядра, можно отобразить данные на многомерном пространстве и, возможно, добиться лучшего соответствия модели данным.

Однако существует множество функций ядра, и бывает сложно выбрать одну из них. В соответствии с распространенным соглашением следует начинать с гауссова *RBF-ядра*, которое, как было показано, хорошо работает для многих типов данных. Для того чтобы обучить SVM-модель на основе RBF-ядра, можно воспользоваться функцией `ksvm()`:

```
> letter_classifier_rbf <- ksvm(letter ~ .,
data =
letters_train,                                ker
nel = "rbfdot")
```

Затем делаем прогнозы, как и раньше:

```
> letter_predictions_rbf <-
predict(letter_classifier_rbf,
letters_test)
```

Наконец, сравним точность с линейной SVM-моделью:

```
> agreement_rbf <- letter_predictions_rbf ==
letters_test$letter>
table(agreement_rbf) agreement_rbf FALSE TRUE 27
```



```
5 3725>
```

```
prop.table(table(agreement_rbf)) agreement_rbf  FA  
LSE      TRUE 0.06875 0.93125
```



Из-за случайности в RBF-ядре `ksvm` ваши результаты могут отличаться от представленных здесь. Если вы хотите, чтобы они точно совпадали, перед запуском функции `ksvm()` используйте команду `set.seed(12345)`.

Всего лишь изменив функцию ядра, мы смогли повысить точность модели распознавания символов с 84 до 93 %.

Определение наилучшего параметра штрафа для SVM

Если этот уровень эффективности все еще недостаточен для OCR-программы, безусловно, можно проверить другие ядра. Однако есть один плодотворный подход. Он заключается в изменении параметра штрафа, от которого зависит ширина границы в SVM-решении. Этот параметр определяет баланс модели между перетренированностью и недотренированностью: чем выше штраф, тем настойчивее алгоритм обучения будет пытаться идеально классифицировать каждый обучающий экземпляр, поскольку за каждую ошибку назначается более высокий штраф. С одной стороны, высокий штраф способен привести к тому, что алгоритм обучения станет перетренированным. С другой стороны, слишком низкий штраф может привести к тому, что алгоритм упустит важные, неявные закономерности в тренировочных данных и окажется недообученным.

Не существует правила, которое позволило бы заранее узнать идеальный размер штрафа, поэтому рассмотрим, как модель работает для различных значений параметра штрафа `C`. Вместо того, чтобы несколько раз повторять процесс обучения и оценки, воспользуемся функцией `sapply()`, которая позволяет применить функцию к вектору значений потенциального штрафа. Сначала мы с помощью функции `seq()` сгенерируем этот вектор как последовательность чисел от 5 до 40 с шагом 5. Затем, как показано в коде далее, пользовательская функция обучает модель, как и раньше, каждый раз используя новое значение штрафа и делая прогнозы для тестового набора данных. Точность каждой модели вычисляется как число прогнозов, которые соответствуют фактическим значениям, разделенное на общее количество прогнозов (рис. 7.22). Результат визуализируется с помощью функции `plot():>`

```
cost_values <- c(1, seq(from = 5, to = 40, by =  
5))>> accuracy_values <- sapply(cost_values,  
function(x) { set.seed(12345) m <-
```

```
ksvm(letter ~ ., data =
letters_train,          kernel = "rbfdot", C
= x)  pred <- predict(m, letters_test)  agree
<- ifelse(pred == letters_test$letter, 1,
0)  accuracy <- sum(agree) /
nrow(letters_test)  return (accuracy) })>
plot(cost_values, accuracy_values, type = "b")
```

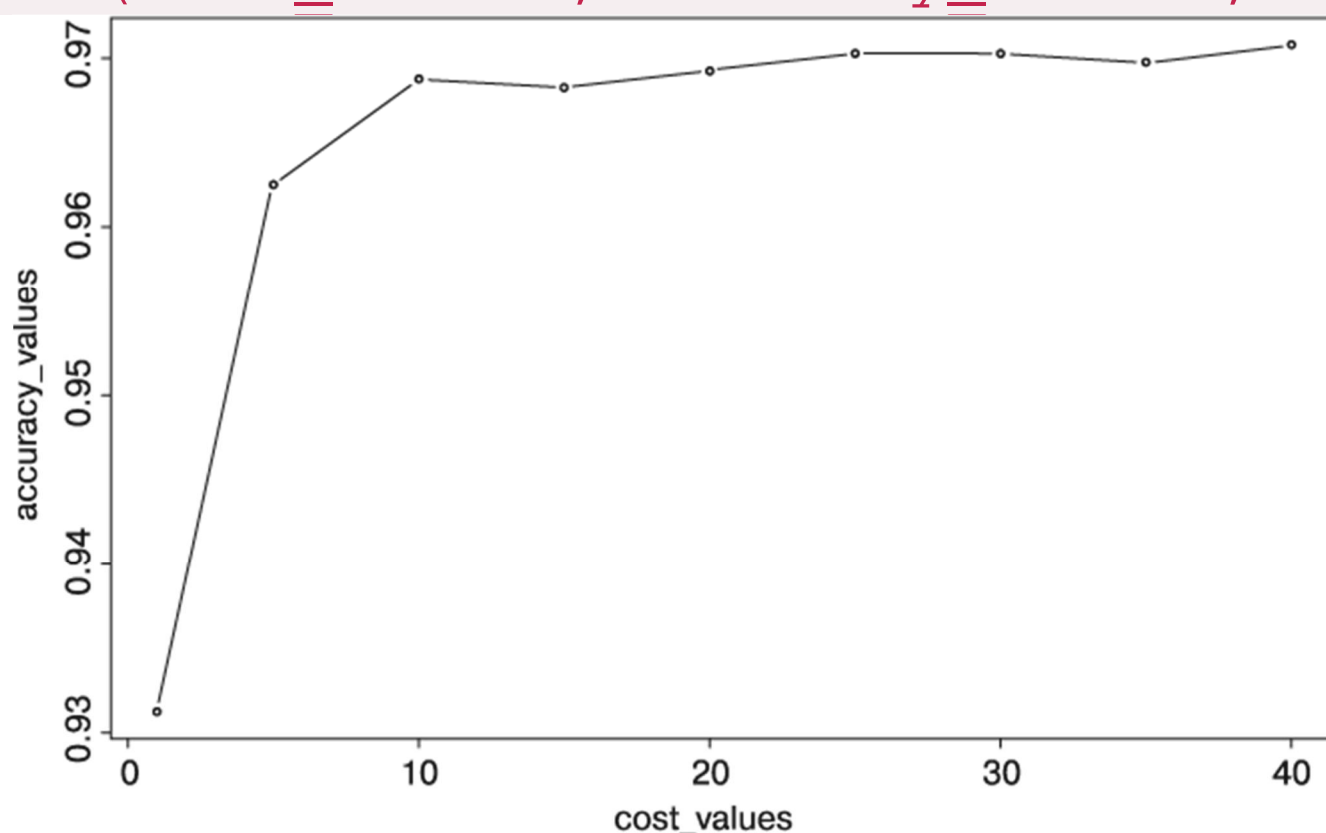


Рис. 7.22. Зависимость точности SVM-алгоритма от штрафа

Как видно на графике, предлагаемое по умолчанию значение штрафа $C=1$ для SVM-алгоритма привело к получению наименее точной из девяти оцененных моделей — ее точность равна 93 %. Если же присвоить C значение 10 или выше, то получим точность около 97 %, что является значительным повышением эффективности! Вероятно, такая модель уже достаточно близка к тому, чтобы развернуть ее в реальной среде, хотя, возможно, стоит поэкспериментировать с различными ядрами, чтобы увидеть, не получится ли приблизиться к 100%-ной точности. Каждое новое повышение точности приводит к уменьшению количества ошибок в OCR-программе и улучшает общее восприятие программы конечным пользователем.

Резюме

В этой главе мы рассмотрели два метода ML, которые имеют большой потенциал, но часто упускаются из виду из-за их сложности. Надеюсь, теперь вы убедились, что эта репутация несколько незаслуженна. Основные понятия, которые лежат в основе нейронных сетей и SVM-алгоритмов, вполне легко усвоить.

С другой стороны, поскольку нейронные сети и SVM-методы существуют уже много десятков лет, у каждого из них появилось

множество модификаций. В этой главе мы лишь в общих чертах рассказали о том, что можно делать с помощью таких методов. Изучив терминологию, вы теперь сможете видеть нюансы, отличающие многие достижения, которые появляются каждый день, включая постоянно растущую область применения методов глубокого обучения.

После изучения моделей прогнозирования, от простых до сложных, в следующей главе рассмотрим методы, применяемые для решения других типов задач обучения. Это неадаптивные методы обучения, которые позволяют выявлять интересные зависимости между данными.

8. Обнаружение закономерностей: анализ потребительской корзины с помощью ассоциативных правил

Вспомните свою последнюю спонтанную покупку. Возможно, это была жевательная резинка или шоколадка на кассе в продуктовом магазине. Или во время ночной поездки за подгузниками и детской смесью вы прихватили напиток с кофеином или пару бутылок пива. Возможно, вы даже эту книгу купили по рекомендации продавца. Все эти импульсивные покупки неслучайны, так как розничные продавцы используют сложные методы поиска закономерностей, чтобы идентифицировать паттерны потребительского поведения при розничной торговле.

В прошлые годы такие рекомендации основывались на субъективной интуиции специалистов по маркетингу и менеджеров службы снабжения. Теперь сканеры штрихкодов, базы данных складов и интернет-магазины накопили достаточно данных о транзакциях, которые можно использовать при машинном обучении для изучения паттернов поведения потребителей. Такая практика широко известна как *анализ потребительской корзины*, так как она часто применяется для анализа данных о супермаркетах.

Несмотря на то что эта технология была основана на анализе покупок, она также полезна и в других областях. Когда вы дочитаете эту главу, сможете применять методы анализа потребительской корзины для своих задач, какими бы они ни были. Как правило, работа этих методов включает в себя следующее:

- использование простых показателей эффективности для поиска взаимосвязей в больших базах данных;
- понимание особенностей данных о транзакциях;
- умение выявлять полезные и действенные закономерности.

Цель анализа потребительской корзины — выявить действенные паттерны. Таким образом, по мере того, как мы будем применять эту методику, вы, скорее всего, сможете определить, как именно ее можно будет

использовать в вашей работе, даже если она не связана с розничными продажами.

Ассоциативные правила

Составными частями анализа потребительской корзины являются товары, которые есть в любой транзакции. Если группа из одного или нескольких элементов заключена в фигурные скобки, это означает, что данные элементы образуют множество — множество товаров, которое встречается в данных с некоторой регулярностью. Транзакции определяются как множества товаров. Ниже приведен пример транзакции, которая встречается в типичном продуктовом магазине:

{хлеб, арахисовое масло, мармелад}.

Результатом анализа потребительской корзины является набор *ассоциативных правил*, которые определяют закономерности, обнаруженные во взаимосвязях между содержимым множества элементов. Ассоциативные правила всегда состоят из двух подмножеств множества элементов и обозначаются как отношение между одним подмножеством элементов, расположенным слева (Left-Hand Side, LHS), и другим подмножеством, расположенным справа (Right-Hand Side, RHS). LHS — это условие, которое должно быть выполнено для запуска правила, а RHS — ожидаемый результат выполнения этого условия. Правило, сформулированное на основе транзакции из предыдущего примера, может быть выражено следующим образом:

{арахисовое масло, мармелад} → {хлеб}.

Проще говоря, это ассоциативное правило гласит, что если покупатель приобретает арахисовое масло и мармелад, то, скорее всего, он также купит хлеб. Другими словами, «покупка арахисового масла и мармелада подразумевает покупку хлеба».

Ассоциативные правила, разработанные на основании информации из баз данных розничных транзакций, используются не для прогнозирования, а для неконтролируемого поиска знаний в больших базах данных, и этим они отличаются от алгоритмов классификации и числового прогнозирования, представленных в предыдущих главах. Тем не менее, как вы вскоре обнаружите, алгоритмы обучения ассоциативных правил тесно связаны и имеют те же особенности обучения, что и методы, изучающие правила классификации, описанные в главе 5.

Поскольку алгоритмы обучения ассоциативных правил являются неконтролируемыми, нет необходимости в обучении алгоритма; данные не должны быть снабжены метками заранее. Программа просто выполняется для набора данных в надежде, что удастся выявить интересные зависимости. Конечно, недостатком этого процесса является то, что не существует простого способа объективно измерить эффективность

алгоритма обучения ассоциативных правил, кроме качественной оценки полезности — как правило, это оценка на глаз в той или иной форме.

Несмотря на то что ассоциативные правила чаще всего используются для анализа потребительской корзины, они полезны при поиске паттернов для многих других типов данных. Другие потенциальные области применения этой технологии включают:

- поиск интересных и часто встречающихся паттернов в последовательностях ДНК и белков среди данных о раковых заболеваниях;
- поиск паттернов покупок или медицинских жалоб в сочетании с мошенническим использованием кредитных карт или страховки;
- выявление вариантов поведения, которые предшествуют отказу покупателей от услуг сотовой связи или модернизации ими систем кабельного телевидения.

Анализ ассоциативных правил используется для поиска интересных зависимостей среди очень большого числа элементов. Люди способны интуитивно понимать такие вещи, но для этого часто требуется быть экспертом или иметь большой опыт, в то время как алгоритм обучения ассоциативных правил способен справиться с задачей за несколько минут или даже секунд. Кроме того, некоторые наборы данных слишком велики и сложны и для человека работа с ними подобна поиску иголки в стоге сена.

Алгоритм Apriori для поиска ассоциативных правил

Большие наборы данных о транзакциях создают проблемы не только для людей, но и для машин. Наборы данных о транзакциях могут быть большими по количеству как транзакций, так и записанных в них элементов или признаков. Проблема в том, что количество потенциальных множеств элементов экспоненциально возрастает по мере увеличения числа признаков. Например, если в наборе могут встретиться или не встретиться k элементов, то существует 2^k возможных множеств элементов, которые могут быть потенциальными правилами. Если продавец продает всего 100 различных товаров, то алгоритм должен изучить порядка $2^{100} = 1,27e + 30$ возможных множеств товаров — казалось бы, невыполнимая задача.

Вместо того чтобы исследовать все эти множества товаров один за другим, более умный алгоритм поиска правил ориентируется на тот факт, что в действительности многие из возможных комбинаций элементов редко встречаются на практике — если встречаются вообще. Например, даже если в магазине продаются товары для автомобилей и женская косметика, то множество {моторное масло, губная помада}, скорее всего, будет встречаться очень редко. Игнорируя эти редкие (и, возможно, менее важные) комбинации, можно сузить область поиска правил до более управляемого размера.

Была проделана большая работа по выявлению эвристических алгоритмов для сокращения количества множеств элементов для поиска. Возможно, наиболее широко используемый подход для эффективного поиска в больших базах данных правил — это Apriori. Созданный в 1994 году Ракешем Агравалом (Rakesh Agrawal) и Рамакришнаном Шрикантом (Ramakrishnan Srikant), алгоритм Apriori стал в некоторой степени синонимом обучения ассоциативных правил. Название алгоритма происходит от того факта, что в нем используется простое априорное (от латинского *a priori*) суждение о свойствах часто встречающихся множеств элементов.

Прежде чем более подробно обсудить этот алгоритм, отмечу, что он, как и все алгоритмы обучения, имеет свои преимущества и недостатки. Некоторые из них перечислены в табл. 8.1.

Таблица 8.1

Преимущества	Недостатки
<p>Способен работать с большими объемами данных о транзакциях.</p> <p>Результаты представлены в виде понятных правил.</p> <p>Полезен для интеллектуального анализа данных и обнаружения неожиданных зависимостей в базах данных</p>	<p>Не очень полезен для небольших наборов данных.</p> <p>Требует усилий, чтобы отделить истинное знание от здравого смысла.</p> <p>Легко сделать ложные выводы из случайных моделей</p>

Как уже отмечалось, в алгоритме Apriori в качестве руководства для сокращения пространства поиска ассоциативных правил используется простое априорное суждение: все подмножества часто встречающегося множества элементов также должны часто встречаться. Такая эвристика известна как *свойство Apriori*. Используя это проницательное наблюдение, можно резко ограничить количество правил, среди которых проводится поиск.

Например, множество {моторное масло, губная помада} может встречаться часто только в случае, если множества {моторное масло} и {губная помада} также часто встречаются. Следовательно, если моторное масло или губную помаду покупают редко, то любое множество, содержащее эти элементы, можно исключить из поиска.



Подробнее об алгоритме Apriori читайте в статье: Agrawal. R., Srikant R. Fast Algorithms for Mining Association Rules // Proceedings of the 20th International Conference on Very Large Databases, 1994. P. 487–499.

Чтобы увидеть, как этот принцип применяется в более реалистичных условиях, рассмотрим простую базу данных о транзакциях. На рис. 8.1 показаны пять транзакций, совершенных в воображаемом магазине подарков при больнице.

ID транзакции	Множества покупок
1	{цветы, открытка с пожеланием выздоровления}
2	{плюшевый мишка, цветы, шарик, шоколадка}
3	{открытка с пожеланием выздоровления, шоколадка, цветы}
4	{плюшевый мишка, шарик, газированная вода}
5	{цветы, открытка с пожеланием выздоровления, газированная вода}

Рис. 8.1. Множества элементов, соответствующие пяти транзакциям в воображаемом магазине подарков при больнице

Глядя на множества покупок, можно прийти к выводу, что существует несколько типичных паттернов покупки. Человек, посещающий больного друга или члена семьи, как правило, покупает открытку с пожеланиями выздоровления и цветы, в то время как для молодых матерей обычно покупают плюшевых мишек и воздушные шары. Такие паттерны заслуживают внимания, потому что встречаются достаточно часто, чтобы вызвать интерес; подключим логику и опыт в данной области, чтобы объяснить эти правила.

Аналогичным образом алгоритм Apriori использует статистические показатели интересности множества элементов для выявления ассоциативных правил в гораздо более крупных базах данных о транзакциях. Из следующих разделов вы узнаете, каким образом алгоритм Apriori определяет эти показатели интересности и каким образом они в сочетании со свойством априори позволяют сократить количество изучаемых правил.

Измерение интересности правила: поддержка и доверие

Считать ли ассоциативное правило интересным, алгоритм определяет по двум статистическим показателям: поддержке и доверию. Устанавливая минимальные пороговые значения для каждого из этих показателей и применяя принцип Apriori, можно существенно снизить количество правил, рассматриваемых алгоритмом, — возможно, даже до того уровня, на котором будут рассматриваться только правила, очевидные с точки зрения здравого смысла. Важно тщательно изучить типы правил, которые будут исключены в соответствии с этими критериями.

Поддержка (support) множества элементов или правила — это частота появления этого правила среди данных. Например, значение поддержки для множества элементов {открытка с пожеланием выздоровления, цветы} в данных больничного магазина подарков равно $3 / 5 = 0,6$. Следовательно, поддержка правила {открытка с пожеланием выздоровления} → {цветы} также составляет 0,6. Поддержку можно вычислить для любого множества элементов или даже для одного элемента; например, поддержка элемента

{шоколадка} составляет $2 / 5 = 0,4$, так как шоколадки встречаются в 40 % покупок. Функция, вычисляющая поддержку для множества элементов X , может быть определена следующим образом:

$$\text{support}(X) = \frac{|\{X \in N\}|}{\text{count}(X)}.$$

Здесь N — количество транзакций в базе данных, а $\text{count}(X)$ — число транзакций, содержащих множество элементов X .

Доверие (confidence) правила — это мера его способности прогнозирования, или точности. Доверие определяется как поддержка множества элементов, содержащего элементы X и Y , разделенная на поддержку множества элементов, содержащего только X :

$$\text{confidence}(X \rightarrow Y) = \frac{\text{support}(X, Y)}{\text{support}(X)}.$$

В сущности, доверие — это доля транзакций, наличие в которых элемента или множества элементов X приводит к наличию в них элемента или множества элементов Y . Следует учитывать, что доверие того, что наличие в покупке X приводит к наличию Y , не равно доверию того, что наличие Y приводит к наличию X . Например, доверие правила {цветы} → {открытка с пожеланиями выздоровления} составляет $0,6 / 0,8 = 0,75$, а доверие правила {открытка с пожеланиями выздоровления} → {цветы} равна $0,6 / 0,6 = 1,0$. Это означает, что при покупке цветов открытку покупают в 75 % случаев, а при покупке открытки цветы покупают в 100 % случаев. Эта информация может быть весьма полезной для руководства магазина подарков.



Возможно, вы обратили внимание на сходство между поддержкой, доверием и правилами вычисления вероятности в байесовском алгоритме, описанными в главе 4. В сущности, $\text{support}(A, B)$ — это то же самое, что $P(A \cap B)$, а $\text{confidence}(A \rightarrow B)$ — то же, что $P(B|A)$. Различия лишь в контексте.

Такие правила, как {открытка с пожеланиями выздоровления} → {цветы}, называются *сильными правилами*, поскольку имеют высокую поддержку и доверие. Один из способов найти наиболее сильные правила — изучить все возможные сочетания товаров в сувенирном магазине, измерить их поддержку и доверие и передать алгоритму только те правила, которые соответствуют определенным интересам. Однако, как уже отмечалось, такая стратегия, как правило, осуществима лишь для самых маленьких наборов данных.

В следующем разделе вы увидите, как алгоритм Apriori использует минимальные уровни поддержки и доверия вместе с принципом Apriori для быстрого поиска сильных правил путем сокращения количества правил до более управляемого уровня.

Построение набора правил по принципу Apriori

Напомним, что принцип Apriori гласит: все подмножества часто встречающегося множества элементов также должны быть частыми. Другими словами, если $\{A, B\}$ является частым, то $\{A\}$ и $\{B\}$ также должны быть частыми. По определению поддержка показывает, как часто множество элементов встречается в данных. Следовательно, если мы знаем, что $\{A\}$ не соответствует желаемому пороговому значению поддержки, то нет оснований рассматривать $\{A, B\}$ или любое другое множество элементов, включающее в себя $\{A\}$; такое множество не может быть частым.

Алгоритм Apriori использует эту логику, чтобы исключить потенциальные ассоциативные правила, не выполняя их фактической оценки. Процесс формирования правил происходит в два этапа.

1. Определение всех множеств товаров, которые соответствуют минимальному порогу поддержки.

2. Создание из этих множеств элементов правил, соответствующих минимальному порогу доверия.

Первый этап выполняется за несколько итераций. Каждая следующая итерация включает в себя оценку поддержки все более крупных множеств элементов. Например, на первой итерации выполняется оценка множеств элементов, состоящих из одного элемента (одноэлементные множества), на второй — множеств из двух элементов, и т.д. Результатом каждой итерации i является множество, состоящее из всех i -элементных множеств, соответствующих минимальному порогу поддержки.

Все множества элементов, полученные на итерации i , объединяются в сочетания, чтобы построить подходящие множества элементов для оценки на итерации $i + 1$. Но принцип Apriori позволяет исключить некоторые из них еще до начала следующей итерации. Так, если $\{A\}$, $\{B\}$ и $\{C\}$ встречаются на первой итерации часто, а $\{D\}$ — нечасто, то на второй итерации будут рассматриваться только $\{A, B\}$, $\{A, C\}$ и $\{B, C\}$. Таким образом, алгоритму остается оценить только три множества элементов, а не шесть, которые были бы сформированы, если бы множества, содержащие D , не были исключены **априори**.

Предположим, что на второй итерации обнаружилось, что $\{A, B\}$ и $\{B, C\}$ встречаются часто, а $\{A, C\}$ — нет. Тогда, несмотря на то что третья итерация обычно начинается с оценки поддержки $\{A, B, C\}$, этот шаг уже не нужен. Почему? Потому что согласно принципу Apriori $\{A, B, C\}$ не может быть частым, поскольку подмножество $\{A, C\}$ не является частым. Следовательно, алгоритм может остановиться, не создавая новые множества элементов на третьей итерации (рис. 8.2).

Итерация	Должен оценить	Частые множества элементов	Нечастые множества элементов
1	{A}, {B}, {C}, {D}	{A}, {B}, {C}	{D}
2	{A, B}, {A, C}, {B, C} {A, D}, {B, D}, {C, D}	{A, B}, {B, C}	{A, C}
3	<u>{A, B, C}</u>		
4	<u>{A, B, C, D}</u>		

Рис. 8.2. Алгоритм Apriori оценивает лишь 7 из 12 возможных множеств элементов

После этого можно переходить ко второму этапу алгоритма Apriori — генерации ассоциативных правил для всех возможных подмножеств данного множества элементов. Например, {A, B} приведет к созданию правил-кандидатов для $\{A\} \rightarrow \{B\}$ и $\{B\} \rightarrow \{A\}$. Эти правила оцениваются по минимальному порогу доверия, и все правила, которые не соответствуют желаемому уровню доверия, исключаются.

Пример: выявление часто покупаемых продуктов в соответствии с ассоциативными правилами

Как отмечено во введении к этой главе, анализ потребительской корзины используется рекомендательными системами, используемыми во многих обычных и интернет-магазинах. Выявленные ассоциативные правила указывают на сочетания товаров, которые часто покупаются вместе. Знание этих паттернов позволяет создать новые способы оптимизации товаров в сети продуктовых магазинов, рекламных акций или раскладки товаров в магазине. Например, если покупатели часто приобретают на завтрак кофе или апельсиновый сок вместе с выпечкой, то, возможно, удастся повысить прибыль, если разместить выпечку поближе к кофе и сокам.

В этом примере мы выполним анализ потребительской корзины на основе данных о транзакциях продуктового магазина. Однако эти методы можно применять ко многим другим типам задач, от рекомендаций фильмов до сайтов знакомств и для обнаружения опасных зависимостей между лекарствами. При этом мы увидим, как алгоритм Apriori способен эффективно обрабатывать потенциально большие наборы ассоциативных правил.

Шаг 1. Сбор данных

В нашем анализе потребительской корзины будут использованы данные о покупках за один месяц работы реального продуктового магазина. Эти данные содержат 9835 транзакций, то есть примерно 327 транзакций в день (примерно 30 транзакций в час в течение 12-часового рабочего дня), что говорит о том, что этот магазин не очень большой и не маленький.



Используемый здесь набор данных был создан на основе набора данных Groceries из R-пакета arules. Для получения дополнительной информации читайте статью: Hahsler M., Hornik K., Reutterer T. Implications of Probabilistic Data Modeling for Mining Association Rules, 2005, а также Gaul W., Vichi M., Weihs C. From Data and Information Analysis to Knowledge Engineering. Studies in Classification, Data Analysis and Knowledge Organization, 2006. P. 598–605.

В типичном продуктовом магазине продается множество разнообразных товаров. Там может быть пять марок молока, дюжина видов стирального порошка и три наименования кофе. Учитывая средний размер магазина в данном примере, предположим, что здесь не очень обеспокоены поиском правил, применимых только к конкретной марке молока или моющего средства. С учетом этого мы удалили торговые марки из данных о покупках. Это позволило сократить количество продуктов до чуть более чем 169 видов, которыми можно управлять, используя широкие категории, такие как куриное мясо, замороженные продукты, маргарин и газированная вода.



Если вы хотите определить более конкретные ассоциативные правила — например, клиенты предпочитают покупать с арахисовым маслом виноградный или клубничный мармелад, — то вам понадобится огромное количество данных о транзакциях. Крупные розничные сети, чтобы найти ассоциации между конкретными брендами, цветами или запахами товаров, используют базы данных о миллионах транзакций.

Есть ли у вас какие-либо предположения о том, какие типы товаров обычно покупают вместе? Является ли обычным сочетанием вино и сыр? Хлеб и масло? Чай и мед? Тщательнее изучим эти данные и посмотрим, подтвердятся ли предположения.

Шаг 2. Исследование и подготовка данных

Данные о транзакциях хранятся в несколько ином формате, чем тот, что мы использовали ранее. В большинстве предыдущих примеров анализа данные были представлены в виде матрицы, строки которой соответствовали примерам элементов, а столбцы — признакам. В матричном формате все примеры должны иметь одинаковый набор признаков.

Однако данные о транзакциях представлены в более свободной форме. Как обычно, каждая строка данных соответствует одному примеру — в данном случае транзакции. Но вместо фиксированного количества

признаков каждая запись содержит разделенный запятыми список, состоящий из произвольного количества элементов, одного или нескольких. По сути, у каждого примера могут быть свои признаки.



Для того чтобы продолжить выполнение этой задачи анализа, загрузите файл `groceries.csv` и сохраните его в рабочем каталоге `R`.

Первые пять строк исходного файла `groceries.csv` выглядят так:

```
citrus fruit,semi-finished
bread,margarine,ready souptropical
fruit,yogurt,coffeewhole milkpip
fruit,yogurt,cream cheese,meat spreadsother
vegetables,whole milk,condensed milk,long life
bakery product
```

Эти строки соответствуют пяти отдельным транзакциям по покупке продуктов. Первая транзакция включает в себя четыре товара: цитрусовые, тесто, маргарин и готовые супы. Для сравнения, третья транзакция состоит всего из одного товара: цельного молока.

Предположим, мы попытались загрузить данные с помощью функции `read.csv()`, как в предыдущих задачах анализа. На рис. 8.3 показано, как `R` представит данные в матричном формате.

Как видим, `R` создал четыре столбца для хранения элементов данных о транзакциях: `V1`, `V2`, `V3` и `V4`. Возможно, покажется целесообразным использовать данные в такой форме, однако если так поступить, то позже возникнут проблемы. `R` решил создать четыре переменные, потому что в первой строке было ровно четыре значения, разделенных запятыми. Однако, как мы знаем, в состав корзины может входить более четырех элементов; в схеме с четырьмя столбцами такие транзакции будут разбиты на несколько строк матрицы. Мы могли бы попытаться исправить это, поместив в начало файла транзакцию с наибольшим количеством элементов, но тогда игнорируется еще одна, более серьезная проблема.

	V1	V2	V3	V4
1	citrus fruit	semi-finished bread	margarine	ready soups
2	tropical fruit	yogurt	coffee	
3	whole milk			
4	pip fruit	yogurt	cream cheese	meat spreads
5	other vegetables	whole milk	condensed milk	long life bakery product

Рис. 8.3. Данные о транзакциях, неправильно загруженные в матричном формате

Структурируя данные таким способом, `R` создал набор признаков, которые хранят не только элементы транзакций, но и их последовательность. Если представить наш алгоритм обучения как попытку найти зависимости между `V1`, `V2`, `V3` и `V4`, то цельное молоко в

признаке `V1` будет означать не то же самое, что цельное молоко в признаке `V2`. Нам нужен набор данных, в котором бы транзакция представлялась не как множество позиций, заполняемых (или не заполняемых) конкретными элементами, а скорее как потребительская корзина, которая содержит или не содержит каждый конкретный товар.

Подготовка данных: создание разреженной матрицы данных о транзакциях

Решение этой проблемы заключается в использовании структуры данных, называемой разреженной матрицей. Возможно, вы помните, что мы использовали разреженную матрицу для обработки текстовых данных в главе 4. Как и в предыдущем наборе данных, каждая строка разреженной матрицы соответствует одной транзакции. Однако в разреженной матрице каждому товару, который может в итоге оказаться в корзине покупателя, соответствует отдельный столбец (признак). Поскольку в наших данных продуктового магазина присутствует 169 различных товаров, разреженная матрица будет содержать 169 столбцов.

Почему бы просто не сохранить это как фрейм данных, как в большинстве предыдущих задач анализа? Причина заключается в том, что при добавлении новых транзакций и товаров обычная структура данных быстро становится слишком большой, и перестает помещаться в памяти компьютера. Даже для относительно небольшого набора данных о транзакциях матрица содержит почти 1,7 миллиона ячеек, в большинстве из которых записаны нули (отсюда и слово «разреженная» в названии матрицы — в ней очень мало ненулевых значений).

Поскольку нет смысла хранить все эти нули, разреженная матрица практически никогда не хранится в памяти полностью; хранятся только те ячейки, в которых что-то записано. Это позволяет размещать структуру в памяти более эффективно, чем матрицу эквивалентного размера или фрейм данных.

Для того чтобы создать разреженную матричную структуру на основе данных о транзакциях, можно воспользоваться функционалом пакета `arules` (от association rules — «ассоциативные правила»).

Установите и загрузите этот пакет с помощью команд `install.packages("arules")` и `library(arules)`.



Подробнее о пакете `arules` читайте здесь: Hahsler M., Gruen B., Hornik K. `arules` — A Computational Environment for Mining Association Rules and Frequent Item Sets // Journal of Statistical Software, 2005. Vol. 14.

Загружая данные о транзакциях, мы не можем, как раньше, просто использовать функцию `read.csv()`. В пакете `arules` есть

функция `read.transactions()`. Она похожа на `read.csv()`, но ее результатом является разреженная матрица, подходящая для данных о транзакциях. Параметр `sep=","` указывает на то, что элементы входного файла разделяются запятой. Чтобы прочитать данные из `groceries.csv` в разреженную матрицу с именем `groceries`, введите следующую строку:

```
> groceries <-  
read.transactions("groceries.csv", sep = ",")
```

Для того чтобы просмотреть основную информацию о созданной нами матрице продуктов, можно воспользоваться функцией `summary()` для объекта:

```
> summary(groceries)transactions as itemMatrix  
in sparse format with9835 rows  
(elements/itemsets/transactions) and169 columns  
(items) and a density of 0.02609146
```

Первый блок информации в выводе функции (как показано выше) представляет собой сводную статистику о разреженной матрице, которую мы создали. `9835rows` означает количество транзакций, а `169columns` — 169 видов товаров, которые могут появиться в корзине покупателя. Каждая ячейка матрицы равна 1, если товар был приобретен в соответствующей транзакции, или 0 в противном случае.

Значение *плотности*, равное `0.02609146` (2,6 %), означает долю ненулевых ячеек матрицы. Поскольку в матрице $9835 \times 169 = 1\,662\,115$ позиций, то, как можно подсчитать, за 30 дней работы магазина было приобретено в общей сложности $1\,662\,115 \times 0,026\,091\,46 = 43\,367$ товаров (не учитывая, что многие товары покупались по нескольку раз). Выполнив еще один шаг, мы можем определить, что средняя транзакция содержала $43\,367 / 9835 = 4409$ различных товаров. Разумеется, если внимательнее посмотреть на результат, то станет ясно, что среднее количество элементов в транзакции там уже есть.

В следующем блоке результатов `summary()` перечисляются элементы, которые чаще всего встречаются в данных о транзакциях. Поскольку $2513 / 9835 = 0,2555$, то легко определить, что цельное молоко встречается в 25,6 % транзакций. В список часто покупаемых товаров также входят овощи, рогалики и булочки, газированная вода и йогурт:

most frequent items:	whole milk	other
vegetables	rolls/buns	2513
1903	1809	soda
yogurt	(Other)	1715
1372	34055	

Наконец, мы можем рассмотреть статистические данные о размере транзакций; 2159 транзакций содержали только один элемент, в то время как самая большая транзакция содержала 32 элемента. Первый квартиль и средний размер покупки составили два и три элемента соответственно. Это означает, что 25 % транзакций содержали два или один товар, а примерно половина транзакций — три товара и меньше. Среднее значение 4409 элементов на транзакцию соответствует значению, которое мы вычислили вручную.

```

element (itemset/transaction) length
distribution:sizes      1      2      3      4      5
6      7      8      9     10     11     12 2159   1643   1299
1005   855   645   545   438   350   246   182   117   13
      14     15     16     17     18     19     20     21     22
      23     24    78     77     55     46     29     14     14
9      11     4      6      1     26     27     28     29     32
      1      1      1      3      1     Min.  1st
Qu.    Median    Mean  3rd
Qu.    Max.1.000    2.000    3.000    4.409    6.
000    32.000

```

В пакете `arules` есть несколько полезных функций для исследования данных о транзакциях. Чтобы просмотреть содержимое разреженной матрицы, используйте функцию `inspect()` в сочетании с векторными операторами R. Первые пять транзакций можно посмотреть следующим образом:

```

> inspect(groceries[1:5]) items1 {citrus
fruit, margarine, ready soups, semi-
finished bread}2 {coffee, tropical
fruit, yogurt}3 {whole milk}4 {cream
cheese, meat spreads, pip fruit, yogurt}5
{condensed milk, long life bakery
product, other vegetables, whole milk}

```

Эти транзакции соответствуют представлению исходного CSV-файла. Чтобы исследовать конкретный элемент (столбец данных), используйте матричную запись `[row, col]`. Использование такой нотации в сочетании с функцией `itemFrequency()` позволяет увидеть долю транзакций, в которых содержится указанный элемент. Например, для того чтобы узнать значение поддержки для первых трех элементов в данных о продуктах, можно воспользоваться следующей командой:

```

> itemFrequency(groceries[, 1:3])abrasive
cleaner artif. Sweetener baby

```

```
cosmetics      0.0035587189      0.0032536858      0.0006100661
```

Обратите внимание, что элементы разреженной матрицы сортируются по столбцам в алфавитном порядке. Абразивные чистящие средства (`abrasivecleaner`) и искусственные подсластители (`artif.sweetener`) встречаются примерно в 0,3 % транзакций, а детская косметика (`babycosmetics`) — примерно в 0,06 % транзакций.

Визуализация поддержки элементов: частотные диаграммы элементов

Чтобы визуально представить эту статистику, воспользуемся функцией `itemFrequencyPlot()`. Эта функция строит столбчатую диаграмму, отображающую долю транзакций, содержащих указанные элементы. Поскольку данные о транзакциях содержат очень большое количество элементов, часто приходится ограничивать количество элементов, отображаемых на графике, чтобы получить читаемую диаграмму.

Для того чтобы отобразить на диаграмме элементы, которые присутствуют в минимальном количестве транзакций, можно использовать функцию `itemFrequencyPlot()` с параметром `support`:

```
> itemFrequencyPlot(groceries, support = 0.1)
```

На рис. 8.4 представлена гистограмма, где показаны восемь элементов из набора данных о продуктах, для которых значение поддержки составляет не менее 10 %.

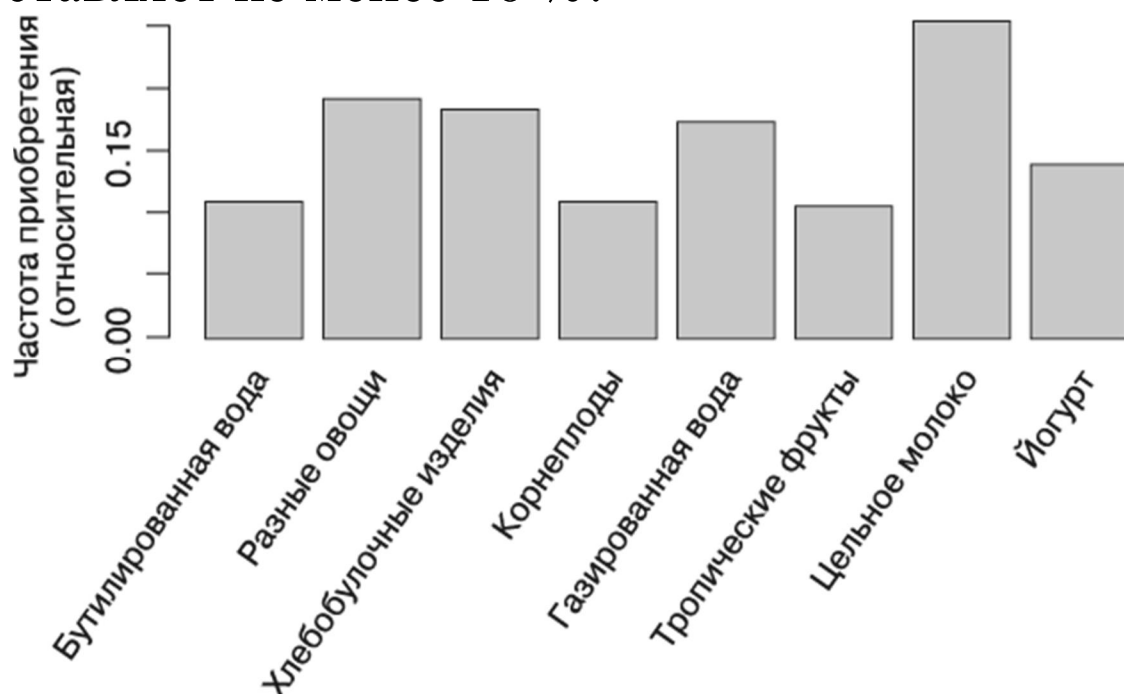


Рис. 8.4. На диаграмме представлены продукты, значение поддержки которых составляет не менее 10 % транзакций

Если вы предпочитаете ограничить диаграмму определенным количеством элементов, используйте `itemFrequencyPlot()` с параметром `topN`:

```
> itemFrequencyPlot(groceries, topN = 20)
```


Затем продукты, представленные на гистограмме, упорядочиваются по убыванию поддержки, как показано на рис. 8.5 для 20 лучших товаров из набора данных о продуктах.

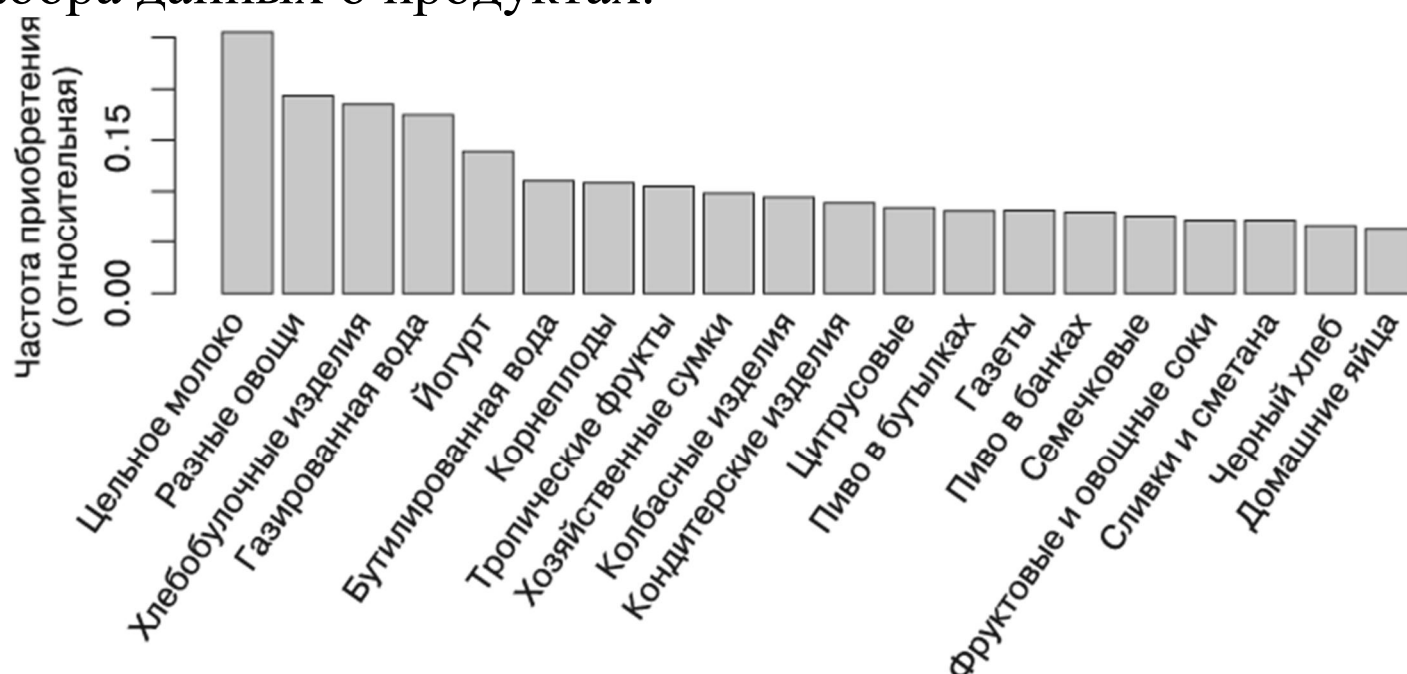


Рис. 8.5. Уровни поддержки для 20 лучших продуктов

Визуализация данных о транзакциях: построение разреженной матрицы

Кроме изучения отдельных элементов, также можно рассмотреть всю разреженную матрицу в целом, используя функцию `image()`. Конечно, поскольку сама матрица слишком велика, лучше запросить какое-то подмножество матрицы. Команда для отображения разреженной матрицы для первых пяти транзакций выглядит следующим образом:

```
> image(groceries[1:5])
```

Как видно на рис. 8.6, диаграмма представляет собой матрицу из пяти строк и 169 столбцов, соответствующих пяти запрошенным нами транзакциям и 169 возможным товарам. Те ячейки матрицы, где в соответствующей транзакции (строке), присутствует данный элемент (столбец), заполняются черным цветом.



Рис. 8.6. Визуализация разреженной матрицы для первых пяти транзакций

Несмотря на то что диаграмма слишком мала и ее, возможно, не очень удобно читать, вы могли заметить, что первая, четвертая и пятая транзакции содержали по четыре элемента, поскольку в их строках заполнены четыре ячейки. Как видно справа на диаграмме, строки 3 и 5, а также строки 2 и 4 имеют общий элемент.

Такая визуализация может быть полезным инструментом для исследования данных о транзакциях. С одной стороны, она может помочь

выявить потенциальные проблемы с данными. Полностью заполненные столбцы могут указывать на товары, которые приобретаются в каждой транзакции, — проблема, которая может возникнуть, например, в случае непреднамеренного включения в набор данных о транзакциях имени или идентификационного номера продавца.

Кроме того, паттерны, выявленные на диаграмме, помогают обнаруживать интересные сегменты транзакций и элементов, особенно если данные отсортированы подходящими способами. Например, если отсортировать транзакции по дате, то паттерны, представленные в виде черных точек, могут указать на сезонные изменения количества или типов приобретаемых товаров. Возможно, на Рождество или Хануку чаще покупают игрушки, а ближе к Хэллоуину спросом пользуются конфеты. Такой тип визуализации может быть особенно полезным, если отсортировать элементы еще и по категориям. Однако в большинстве ситуаций график будет довольно случайным, как статические помехи на телевизионном экране.

Имейте в виду, что для сверхбольших баз данных транзакций такая визуализация не принесет особой пользы, так как ячейки будут слишком маленькими и их невозможно будет различить. Однако в сочетании с функцией `sample()` вы сможете просматривать разреженную матрицу для случайного набора транзакций. Команда для создания случайного набора из 100 транзакций выглядит следующим образом:

```
> image(sample(groceries, 100))
```

В результате будет построена матричная диаграмма со 100 строками и 169 столбцами (рис. 8.7).

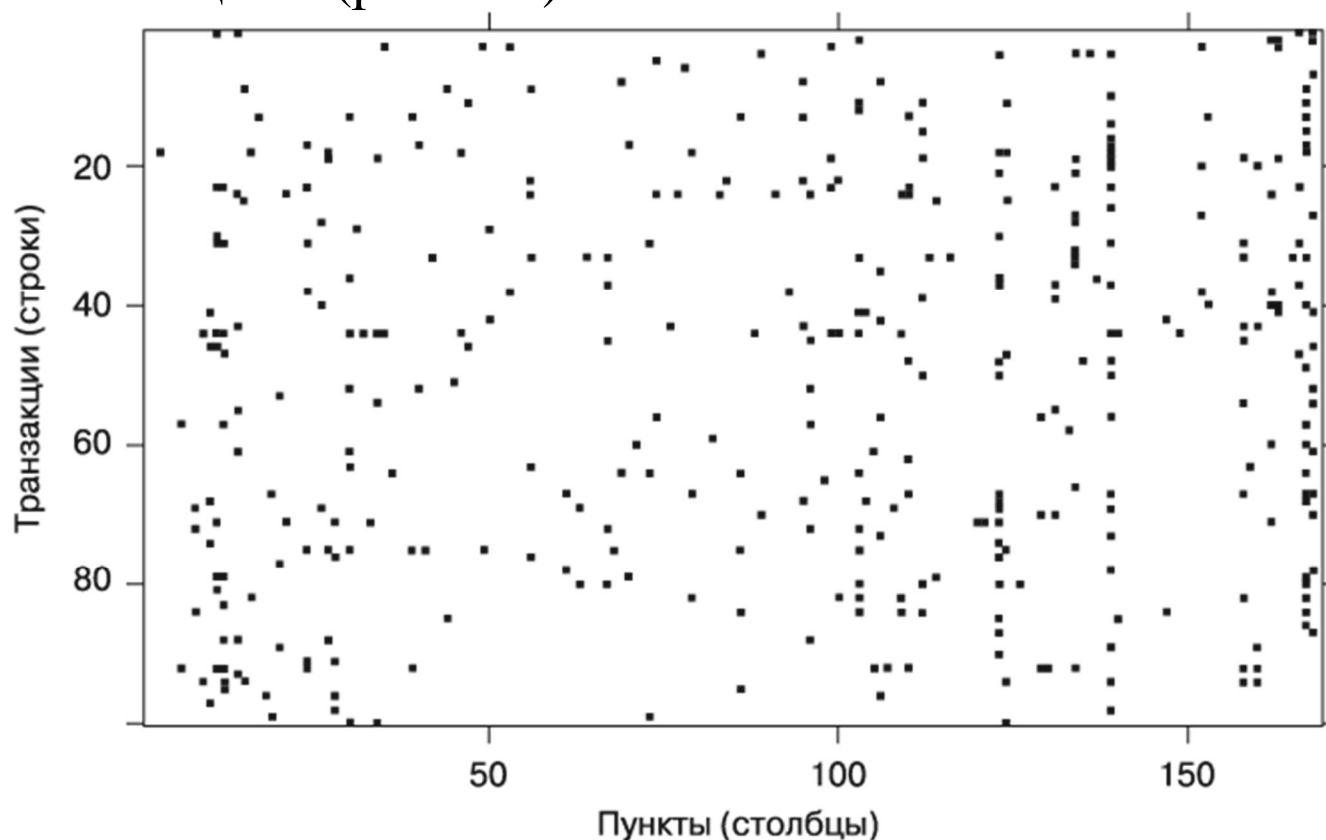


Рис. 8.7. Визуализация разреженной матрицы для 100 случайных транзакций

Несколько столбцов кажутся довольно густо заполненными. Это говорит о наличии в магазине нескольких очень популярных товаров. Тем не менее

в целом распределение точек представляется довольно случайным. Похоже, мы здесь больше ничего не найдем, так что продолжим анализ.

Шаг 3. Обучение модели на данных

Завершив подготовку данных, можно поискать зависимости между товарами в потребительской корзине. Воспользуемся реализацией алгоритма Apriori из того же пакета `arules`, который мы использовали для исследования и подготовки данных о продуктах. Вам нужно установить и загрузить этот пакет, если вы еще этого не сделали. В следующей таблице показан синтаксис для создания наборов правил с помощью функции `apriori()`.

Синтаксис построения ассоциативных правил
Использование функции <code>apriori()</code> из пакета <code>arules</code>
Поиск ассоциативных правил: <pre>myrules <- apriori(data = mydata, parameter = list(support = 0.1, confidence = 0.8, minlen = 1))</pre> <p><code>data</code> – разреженная матрица данных о транзакциях; <code>support</code> – минимально необходимая поддержка правила; <code>confidence</code> – минимально необходимое доверие правила; <code>minlen</code> – минимальное число элементов в транзакции.</p> <p>Функция возвращает объект правил, в котором хранятся все правила, удовлетворяющие указанным минимальным критериям.</p>
Проверка ассоциативных правил: <pre>inspect(myrules)</pre> <p><code>myrules</code> – множество ассоциативных правил, полученных с помощью функции <code>apriori()</code>.</p> <p>Функция выводит ассоциативные правила на экран. Для выбора одного или нескольких правил для просмотра можно применять к <code>myrules</code> векторные операторы.</p>
Пример: <pre>groceryrules <- apriori(groceries, parameter = list(support = 0.01, confidence = 0.25, minlen = 2)) inspect(groceryrules[1:3])</pre>

Функция `apriori()` выполняется просто, но иногда только методом проб и ошибок можно найти значения поддержки и доверия, которые позволят получить разумное количество ассоциативных правил. С одной стороны, если установить эти уровни слишком высокими, то можно не найти правил или найти такие правила, которые окажутся слишком общими, чтобы быть полезными. С другой — слишком низкий порог может привести к громоздкому набору правил. Хуже того, операция может занять очень много времени или исчерпать память на этапе обучения.

Если использовать для наших данных о продуктах предлагаемые по умолчанию значения `support=0.1` и `confidence=0.8`, то получится набор, состоящий из нуля правил:

```
> apriori(groceries)set of 0 rules
```

Очевидно, нам нужно немного расширить поиск.



Если вдуматься, то такой результат не вызывает удивления. Поскольку по умолчанию `support = 0.1`, то для генерации правила элемент должен появляться как минимум в $0,1 \times 9,385 = 938,5$ транзакциях. Поскольку в наших данных так часто встречаются всего восемь элементов, то неудивительно, что мы не нашли подходящих правил.

Один из способов решить проблему выбора минимального значения поддержки состоит в том, чтобы подумать, при каком наименьшем количестве транзакций паттерн будет представлять интерес. Например, можно решить, что если некий товар покупается дважды в день (примерно 60 раз в месяц), то это может быть важно. Таким образом, мы можем вычислить значение поддержки, необходимое для поиска правил, соответствующих как минимум такому количеству транзакций. Поскольку 60 из 9835 транзакций — это 0,006, то мы сначала попробуем установить уровень поддержки таким.

Выбор минимального доверия предполагает поиск тонкого баланса. С одной стороны, если уровень доверия слишком низок, то алгоритм завалит нас большим количеством недостоверных правил — например, десятками правил, указывающих на товары, обычно покупаемые в комплекте с аккумуляторами. Как узнать, на что направить рекламный бюджет? С другой стороны, если установить слишком высокий уровень доверия, то мы будем ограничены правилами, которые очевидны или неизбежны, — например, такими, что в сочетании с аккумуляторами всегда покупается детектор дыма. В этом случае перемещение детекторов дыма ближе к аккумуляторам едва ли принесет дополнительный доход, поскольку эти два предмета и так почти всегда покупают вместе.



Подходящий минимальный уровень доверия во многом зависит от целей анализа. Если начать с консервативного значения, то его всегда можно уменьшить, чтобы расширить поиск, если до сих пор не было найдено ничего полезного.

Начнем с порога доверия, равного 0,25. Это означает, что для того, чтобы правило попало в результаты, оно должно быть корректным по крайней мере в 25 % случаев. Таким образом мы исключим самые ненадежные правила, одновременно оставляя себе некоторое пространство для маневра,

чтобы позже изменить поведение покупателей с помощью целевых рекламных акций.

Теперь мы готовы сгенерировать некоторое количество правил. Кроме минимальных значений поддержки и доверия, полезно установить `minlen=2`, чтобы исключить правила, которые содержат менее двух элементов. Так мы предотвратим создание неинтересных правил, существующих только потому, что какой-то продукт часто покупается, например, `{ }=>wholemilk`. Данное правило выполняет требования минимальной поддержки и доверия, потому что цельное молоко покупается более чем в 25 % транзакций, но это не очень полезная информация.

Полная команда для поиска набора ассоциативных правил с использованием алгоритма Apriori выглядит следующим образом:

```
> groceryrules <- apriori(groceries, parameter = list(support = 0.006, confidence = 0.25, minlen = 2))
```

В результате наши правила будут сохранены в объекте правил, для просмотра которого достаточно ввести его имя:

```
> groceryruleset of 463 rules
```

Наш объект `groceryrules` содержит набор из 463 ассоциативных правил. Для того чтобы определить, являются ли какие-либо из них полезными, нужно копнуть глубже.

Шаг 4. Оценка эффективности модели

Для того чтобы получить общее представление об ассоциативных правилах, можно воспользоваться функцией `summary()` следующим образом. Распределение правил по длине говорит о том, сколько правил состоит из соответствующего количества элементов. В нашем наборе 150 правил состоят из двух элементов, 297 — из трех, а 16 — из четырех. Также в выводе функции представлена сводная статистика, связанная с этим распределением:

```
> summary(groceryrules)set of 463 rulesrule length distribution (lhs + rhs): sizes 2 3 4 150 297 16 Min. 1st Qu. Median Mean 3rd Qu. Max. 2.000 2.000 3.000 2.711 3.000 4.000
```



Как видно в представленном выше выводе функции `summary()`, размер правила равен сумме его левой (`lhs`) и правой частей (`rhs`). Например, размер правила `{хлеб} => {масло}` равен двум элементам, а размер правила `{арахисовое масло, мармелад} => {хлеб}` — трем.

Далее мы видим сводную статистику показателей качества правила: значения поддержки (`support`), доверия (`confidence`) и лифта (`lift`). Параметры `support` и `confidence` не должны вызывать удивления, поскольку мы использовали их в качестве критериев отбора правил. Если бы значения поддержки и доверия для большинства или всех правил оказались очень близкими к минимальным порогам, это, вероятно, вызвало бы беспокойство, так как могло бы означать, что мы, возможно, установили слишком высокую планку. Однако в данном случае это не так, поскольку многие правила имеют гораздо более высокие значения каждого из этих показателей:

```
summary of quality
measures:      support      confidence      lift
iftMin.       :0.006101  Min.         :0.2500  Min.         :0.99
321st Qu.:0.007117  1st Qu.:0.2971  1st
Qu.:1.6229Median :0.008744  Median
:0.3554  Median
:1.9332Mean      :0.011539  Mean        :0.3786  Mean        :
2.03513rd Qu.:0.012303  3rd Qu.:0.4495  3rd
Qu.:2.3565Max.   :0.074835  Max.        :0.6600  Max.
:3.9565
```

Третий столбец — это показатель, который нам еще не встречался. *Лифт* правила показывает, насколько вероятность приобретения одного или нескольких товаров выше, если известно, что при этом был приобретен другой товар или несколько товаров. Лифт вычисляется по следующей формуле:

$$\text{lift}(X \rightarrow Y) = \frac{\text{confidence}(X \rightarrow Y)}{\text{support}(Y)}$$



В отличие от доверия, когда имеет значение последовательность элементов, $\text{lift}(X \rightarrow Y)$ — это то же самое, что $\text{lift}(Y \rightarrow X)$.

Например, предположим, что в продуктовом магазине большинство покупателей приобретают молоко и хлеб. Очевидно, следует ожидать много транзакций, в которых присутствует молоко или хлеб. Однако если $\text{lift}(\text{молоко} \rightarrow \text{хлеб})$ больше единицы, значит, эти два товара вместе встречаются чаще, чем по отдельности. Поэтому большое значение лифта

является явным признаком важности правила и отражает истинную зависимость между элементами.

В последнем фрагменте вывода `summary()` содержится информация об анализе данных, где показано, каким образом были выбраны правила. Здесь мы видим, что данные о продуктах содержали 9835 транзакций и использовались для построения правил с минимальной поддержкой 0,006 и минимальным доверием 0,25:

```
mining
info:      data transactions support confidence
egroceries      9835      0.006      0.25
```

Для того чтобы увидеть конкретные правила, можно воспользоваться функцией `inspect()`. Например, чтобы просмотреть первые три правила в объекте `groceryrules`, можно использовать следующую команду:

```
lhs      rhs      support confidence lift
1 {potted plants} => {whole milk} 0.006914082 0.4000000 1.565460
2 {pasta}      => {whole milk} 0.006100661 0.4054054 1.586614
3 {herbs}      => {root vegetables} 0.007015760 0.4312500 3.956477
```

Первое правило в переводе на обычный язык гласит: «Если покупатель приобретает горшочные растения, то он также покупает цельное молоко». Из значений поддержки около 0,007 и доверия 0,400 следует, что это правило охватывает около 0,7 % транзакций и является верным для 40 % покупок, в состав которых входят растения. Значение лифта говорит о том, насколько более вероятно, что покупатель купит цельное молоко, если он приобрел горшочное растение, по сравнению с обычным покупателем. Поскольку мы знаем, что примерно 25,6 % покупателей купили цельное молоко (`support`) и 40 % покупателей, приобретающих растение, купили цельное молоко (`confidence`), то лифт равен $0,40 / 0,256 = 1,56$, что соответствует значению, показанному в выводе функции `summary()`.



Обратите внимание, что в столбце `support` указывается поддержка для всего правила, а не только для `lhs` или `rhs`.

Несмотря на высокие значения доверия и лифта, считаете ли вы полезным правило {горшочное растение}→{цельное молоко}? Скорее всего, нет, поскольку, похоже, не существует логической причины, по которой некто будет с большей вероятностью покупать молоко, если он купил растение. Тем не менее, наши данные говорят об обратном. Как объяснить этот факт?

Общий подход заключается в том, чтобы разделить ассоциативные правила на следующие категории:

- полезные;
- тривиальные;

необъяснимые.

Очевидно, что цель анализа потребительской корзины — найти **полезные** правила, которые дают четкую и полезную информацию. Не все правила одновременно и понятны, и полезны; сочетание обоих этих факторов встречается гораздо реже, чем каждое из них в отдельности.

Так называемые **тривиальные** правила — это любые правила, которые настолько очевидны, что о них не стоит упоминать: они понятны, но бесполезны. Предположим, что вы являетесь консультантом по маркетингу. Вам платят большие деньги за поиск новых возможностей для совместного продвижения товаров. Если вы сообщите о том, что {детские подгузники} → {молочная смесь}, то вас едва ли пригласят на консультацию еще раз.



Тривиальные правила могут быть замаскированы под более интересные результаты. Предположим, что вы нашли взаимосвязь между определенной маркой детских хлопьев и определенным фильмом на DVD. Такой вывод не очень информативен, если главный герой фильма нарисован на коробке с хлопьями.

Правила являются **необъяснимыми**, если взаимосвязь между элементами настолько неясна, что понять, как использовать эту информацию, невозможно или почти невозможно. Правило может быть просто случайным паттерном в данных — например, правило, утверждающее, что {соленые огурцы} → {шоколадное мороженое}, может быть связано с единственным клиентом, чья беременная жена регулярно испытывала тягу к странным сочетаниям продуктов.

Лучшие правила — это скрытые жемчужины: неизвестные до сих пор идеи, которые становятся очевидными только после того, как были открыты. При наличии достаточного количества времени можно было бы исследовать каждое правило, чтобы найти эти драгоценные камни. Однако исследователи данных, выполняющие анализ, не всегда являются лучшими судьями относительно того, является ли правило полезным, тривиальным или необъяснимым. Поэтому лучшие правила могут появиться благодаря сотрудничеству с экспертами в конкретной предметной области — в данном случае в области управления сетью розничной торговли, — которые способны помочь в интерпретации результатов. В следующем разделе вы узнаете, как организовать такую совместную работу, используя методы сортировки и экспорта изученных правил таким образом, чтобы самые интересные результаты оказались первыми в списке.

Шаг 5. Повышение эффективности модели

Эксперты в данной предметной области очень быстро определяют, какие правила являются полезными, но стало бы неэффективной тратой времени просить их оценить сотни или тысячи правил. Поэтому было бы полезно отсортировать правила по различным критериям и извлечь их из среды R в форме, которую можно передать в отдел маркетинга для более детального изучения. Таким образом можно повысить эффективность правил, сделав результаты более полезными.

Сортировка набора ассоциативных правил

В зависимости от целей анализа потребительской корзины самыми полезными могут оказаться правила с наибольшей поддержкой, доверием или лифтом. В состав пакета `arules` входит функция `sort()`, которая применяется для изменения последовательности в списке правил, позволяя ставить на первое место правила, имеющие самые высокие или самые низкие качественные показатели.

Чтобы изменить последовательность правил в объекте `groceryrules`, можно использовать `sort()`, присвоив параметру `by` значение `"support"`, `"confidence"` или `"lift"`. Сочетая сортировку с векторными операторами, можно получить несколько интересных правил. Например, с помощью следующей команды можно получить пять лучших правил по статистике `lift`:

```
> inspect(sort(groceryrules, by = "lift")[1:5])
```

Результат будет следующий:

lhs	rhs	support	confidence	lift
1 {herbs}	=> {root vegetables}	0.007015760	0.4312500	3.956477
2 {berries}	=> {whipped/sour cream}	0.009049314	0.2721713	3.796886
3 {other vegetables, tropical fruit, whole milk}	=> {root vegetables}	0.007015760	0.4107143	3.768074
4 {beef, other vegetables}	=> {root vegetables}	0.007930859	0.4020619	3.688692
5 {other vegetables, tropical fruit}	=> {pip fruit}	0.009456024	0.2634561	3.482649

Эти правила более интересны, чем те, которые нам встречались ранее. Первое из них, у которого значение `lift` равно примерно 3,96, говорит о том, что люди, которые покупают зелень, почти в четыре раза чаще покупают корнеплоды, чем обычные посетители магазина, — возможно, для приготовления рагу? Правило второе тоже интересно. Взбитые сливки более чем в три раза чаще встречаются в тележке вместе с ягодами, чем в других случаях. Возможно, это какое-то сочетание для десерта?



По умолчанию правила отсортированы в порядке убывания, то есть самые большие значения идут первыми. Чтобы изменить эту последовательность, добавьте в функцию `inspect()` дополнительный параметр `decreasing = FALSE`.

Выделение подмножеств ассоциативных правил

Предположим, что в отделе маркетинга учли предыдущее правило и преисполнились энтузиазмом от открывшихся возможностей создания рекламы для продвижения ягод (сейчас как раз сезон ягод). Однако прежде, чем запускать кампанию, они попросили вас выяснить, как часто ягоды покупаются с другими товарами. Чтобы ответить на этот вопрос, нужно найти все правила, которые включают в себя ягоды в той или иной форме.

Функция `subset()` представляет собой способ поиска подмножеств транзакций, элементов или правил. Чтобы использовать эту функцию для поиска любых правил, в которых встречаются ягоды, нужно воспользоваться следующей командой. В результате правила будут сохранены в новом объекте с именем `berryrules`:

```
> berryrules <- subset(groceryrules, items %in%  
"berries")
```

Затем можно исследовать эти правила, как мы это сделали раньше для всего набора:

```
> inspect(berryrules)
```

В результате получим следующий набор правил:

	lhs	rhs	support	confidence	lift
1	{berries}	=> {whipped/sour cream}	0.009049314	0.2721713	3.796886
2	{berries}	=> {yogurt}	0.010574479	0.3180428	2.279848
3	{berries}	=> {other vegetables}	0.010269446	0.3088685	1.596280
4	{berries}	=> {whole milk}	0.011794611	0.3547401	1.388328

Из четырех правил, касающихся ягод, два кажутся достаточно интересными, чтобы их можно было назвать полезными. Кроме взбитых сливок, ягоды также часто покупают с йогуртом — сочетание, которое хорошо подходит для завтрака или ланча, а также в качестве десерта.

У функции `subset()` много возможностей. Критерии выбора подмножества могут определяться с помощью следующих ключевых слов и операторов.

- Уже встречавшееся нам раньше ключевое слово `items` обозначает элемент, который может появляться в любом месте правила. Чтобы ограничить подмножество правилами, в которых элемент присутствует только с левой или только с правой стороны, вместо этого слова следует использовать `lhs` или `rhs` соответственно.
- Оператор `%in%` означает, что в правиле должен присутствовать хотя бы один из элементов, входящих в указанный список. Если вы хотите отобрать все правила,

где присутствуют ягоды или йогурт, то можете написать `items%in%c("berries", "yogurt")`.

- Также есть дополнительные операторы для частичного (`%pin%`) и полного соответствия (`%ain%`). Частичное соответствие позволяет отобрать правила для цитрусовых и тропических фруктов, используя общий критерий поиска `items%pin%"fruit"`. Для полного соответствия требуется наличие всех перечисленных в списке элементов. Например, по критерию `items%ain%c("berries", "yogurt")` будут найдены правила, в которых присутствуют и ягоды, и йогурт.
- Размеры подмножеств можно ограничить с помощью параметров `support`, `confidence` и `lift`. Например, при `confidence>0.50` будут отобраны только те правила, у которых доверие выше 50 %.
- Критерии отбора можно комбинировать, используя стандартные логические операторы R, такие как `AND(&)`, `OR(|)` и `NOT(!)`.

С помощью этих параметров можно сделать выбор правил настолько конкретным или обобщенным, насколько вам нужно.

Сохранение ассоциативных правил в файле или фрейме данных

Чтобы поделиться результатами анализа потребительской корзины, можно сохранить правила в CSV-файле с помощью функции `write()`. В результате получится CSV-файл, который можно использовать в большинстве редакторов электронных таблиц, включая Microsoft Excel:

```
> write(groceryrules, file =  
"groceryrules.csv", sep = ",", quote =  
TRUE, row.names = FALSE)
```

Иногда также удобно преобразовать правила во фрейм данных R. Это можно сделать с помощью функции `as()` следующим образом:

```
> groceryrules_df <- as(groceryrules,  
"data.frame")
```

В результате будет создан фрейм данных с правилами, представленными в формате фактора, и числовыми векторами для значений поддержки, доверия и лифта:

```
> str(groceryrules_df) 'data.frame':   463 obs.  
of 4 variables:$ rules      : Factor w/ 463 levels  
"{baking powder} =>  
{other          vegetables}",...: 340 302 207  
206 208 341 402 21 139 140 ...$ support    :  
num      0.00691 0.0061 0.00702 0.00773 0.00773  
...$ confidence: num      0.4 0.405 0.431 0.475
```

```
0.475 ...$ lift : num 1.57 1.59 3.96 2.45
1.86 ...
```

Сохранение правил во фрейме данных может быть полезно, если нужно выполнить дополнительную обработку правил или экспортировать их в другую базу данных.

Резюме

Ассоциативные правила применяются для поиска полезной информации в больших базах данных транзакций крупных сетей розничной торговли. Поскольку это обучение без учителя, алгоритмы обучения ассоциативных правил способны извлекать знания из больших баз данных, заранее не имея информации о том, какие закономерности следует искать. Суть этих алгоритмов в том, что они позволяют сократить объем информации до меньшего, более управляемого множества результатов. Алгоритм Apriori, рассмотренный в данной главе, делает это, устанавливая минимальные пороги «интересности» и сообщая только о тех зависимостях, которые удовлетворяют этим критериям.

Мы выполнили анализ по алгоритму Apriori для корзины ежемесячных транзакций небольшого продуктового магазина. Даже в этом простом примере было выявлено множество ассоциаций. Среди них мы выделили несколько паттернов, которые могут быть полезными в будущем для маркетинговых кампаний. Используемые нами методы применимы также для гораздо более крупных сетей розничной торговли с базами данных, многократно превышающими размер рассмотренной базы, а также для проектов, не имеющих отношения к розничным продажам.

В следующей главе изучим другой алгоритм обучения без учителя. Как и ассоциативные правила, он предназначен для поиска закономерностей в данных. Но, в отличие от ассоциативных правил, которые ищут группы взаимосвязанных элементов или признаков, методы, которые будут рассмотрены далее, направлены на поиск зависимостей и взаимосвязей между примерами.

9. Поиск групп данных: кластеризация методом k-средних

Случалось ли вам проводить время, наблюдая за толпой? Если да, то, вероятно, вам встречались типичные персонажи. Возможно, вы видели солидных людей в свежееотглаженных рубашках и с портфелями — в таких типажах легко узнать менеджеров крупных компаний, их еще называют «большими шишками». Или 20-летних парней в джинсах-скинни, фланелевых рубашках и солнцезащитных очках — «хипстеров». Или женщину, которая высаживает детей из минивэна, — «маму-наседку».

Конечно, такие стереотипы опасно применять к людям, поскольку двух одинаковых людей не существует. Тем не менее подобные метки, если рассматривать их как способ описания какого-то собирательного образа, охватывают основные общие черты, присущие представителям данной группы.

Как вы вскоре поймете, процесс кластеризации или обнаружения закономерностей в данных мало отличается от выявления закономерностей среди групп людей. Из этой главы вы узнаете:

- чем задачи кластеризации отличаются от задач классификации, которые мы рассмотрели ранее;
- как кластеризация определяет группу и как такие группы идентифицируются методом *k*-средних — классическим и простым для освоения алгоритмом кластеризации;
- какие шаги необходимо выполнить для применения кластеризации для решения реальной задачи выявления маркетинговых сегментов среди подростков, пользующихся социальными сетями.

Прежде чем приступить к делу, подробно рассмотрим, что подразумевается под кластеризацией.

Что такое кластеризация

Кластеризация — это неконтролируемая задача машинного обучения, которая заключается в автоматическом разделении данных на *кластеры* — группы похожих элементов. Решая эту задачу, алгоритм не знает заранее, как должны выглядеть группы. Поскольку мы не всегда знаем, что именно ищем, кластеризация используется не для прогнозирования, а для обнаружения знаний. Ее цель — найти естественные группировки, существующие в данных.

Каким образом компьютер, не имея глубоких знаний о том, что должно входить в кластер, может знать, где заканчивается одна группа и начинается другая? Ответ прост: в основе кластеризации лежит принцип, гласящий, что элементы внутри кластера должны быть очень похожи между собой, но сильно отличаться от элементов вне этого кластера. Понятие сходства может различаться в разных приложениях, но основная идея всегда одна и та же: сгруппировать данные так, чтобы взаимосвязанные элементы попадали в один кластер.

Полученные кластеры могут быть использованы для различных действий. Например, методы кластеризации применяются в таких областях, как:

- сегментирование потребителей по группам с похожими демографическими характеристиками или покупательскими привычками для организации целевых маркетинговых кампаний;

- обнаружение аномального поведения, такого как несанкционированные вторжения в сеть, путем выявления паттернов использования, выходящих за пределы известных кластеров;
- упрощение чрезвычайно больших наборов данных путем группировки похожих значений признаков в меньшее количество однородных категорий.

Вообще, кластеризация полезна всегда, когда разнообразные, несхожие между собой данные могут быть представлены гораздо меньшим числом групп. Кластеризация позволяет построить осмысленные, полезные структуры данных, которые уменьшают сложность задачи и дают представление о типичных взаимосвязях.

Кластеризация как задача машинного обучения

Кластеризация несколько отличается от рассмотренных нами ранее задач классификации, числового прогнозирования и обнаружения закономерностей. Цель каждого из этих типов задач состоит в том, чтобы построить модель, которая связывает признаки с результатом или одни признаки с другими. Все эти задачи описывают уже существующие паттерны в данных. Целью кластеризации, напротив, является создание новых данных. При кластеризации немаркированным примерам присваивается новая метка кластера, которая логически выводится исключительно на основании взаимосвязей между данными. По этой причине задачу кластеризации иногда называют *неконтролируемой классификацией*, поскольку она в некотором смысле классифицирует немаркированные примеры.

Суть в том, что метки классов, полученные в результате неконтролируемой классификации, не имеют внутреннего значения. Кластеризация покажет, какие группы примеров тесно взаимосвязаны, — например, она может сформировать группы А, В и С, — но только вы сами можете выбрать для этих групп полезные и осмысленные метки. Чтобы увидеть, как это влияет на задачу кластеризации, рассмотрим гипотетический пример.

Предположим, вы организовали конференцию по науке о данных. Чтобы облегчить профессиональное общение и сотрудничество, вы планируете разделить участников на группы в соответствии с одной из трех исследовательских специальностей: информатика, математика и статистика и машинное обучение. К сожалению, разослав приглашения на конференцию, вы поняли, что забыли уточнить, в какой дисциплине специализируется данный участник.

В какой-то момент вы догадались, что можно узнать исследовательскую специализацию каждого ученого, изучив его историю публикаций. Для этого вы собрали данные о количестве статей, опубликованных каждым участником конференции в журналах по информатике, а также по

математике и статистике. Используя данные обо всех участниках конференции, вы построили диаграмму рассеяния (рис. 9.1).



Рис. 9.1. Визуализация ученых по данным их статей по математике и информатике

Как и ожидалось, здесь, похоже, есть закономерность. Можно предположить, что верхний левый угол, в котором представлены те, у кого много публикаций по информатике, но мало статей по математике, могут быть кластером ученых, специализирующихся в области компьютерных наук. Следуя этой логике, в правом нижнем углу находится группа математиков. Аналогичным образом правый верхний угол занимают специалисты по машинному обучению, которые пишут статьи и по математике, и по информатике. Применив эти метки, получим результат, представленный на рис. 9.2.

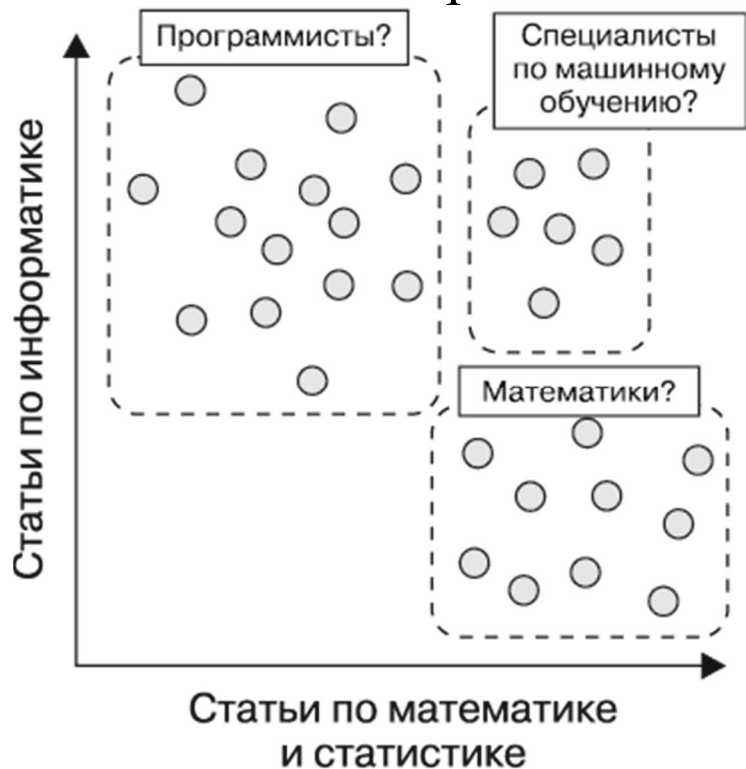


Рис. 9.2. Кластеры, построенные на основе предположений о специализации каждой группы ученых

Мы сформировали группы на глаз — просто определили кластеры как тесно сгруппированные точки данных. Однако, несмотря на кажущиеся очевидными группировки, не спрашивая лично каждого ученого о его специализации, мы не сможем узнать, действительно ли эти группы

однородны. Метки представляют собой предположительные качественные суждения о специализации ученых в каждой группе, основанные на ограниченном наборе количественных данных.

Вместо того чтобы определять границы групп субъективно, было бы неплохо использовать машинное обучение для их объективного определения. Учитывая, что разделение на группы на рис. 9.2 происходит параллельно осям координат, к нашей задаче, судя по всему, применимы деревья решений, описанные в главе 5. Этот метод дал бы нам четкое правило — например, такое: «Если у ученого мало публикаций по математике, то он является экспертом в области информатики». К сожалению, с такой схемой есть сложности. Не имея информации об истинном значении класса для каждой точки данных, контролируемый алгоритм обучения не смог бы сформировать такой паттерн, поскольку у него не было бы способа узнать, какие разделения приведут к однородным группам.

В алгоритмах кластеризации, напротив, применяется процесс, очень похожий на тот, который использовали мы сами, визуально изучая диаграмму рассеяния. Чтобы идентифицировать однородные группы, можно взять за основу меру того, насколько тесно связаны примеры. В следующем разделе разберемся, как реализованы алгоритмы кластеризации.



В этом примере показано интересное применение кластеризации. Если исходные данные являются немаркированными, то с помощью кластеризации можно создать метки классов. Затем можно применить контролируемый алгоритм обучения, такой как деревья решений, чтобы выделить наиболее важные предикторы этих классов. Такой подход называется полуконтролируемым обучением.

Алгоритм кластеризации методом k -средних

Метод k -средних является, пожалуй, наиболее часто используемым алгоритмом кластеризации. Он изучается на протяжении нескольких десятков лет и является основой для многих более сложных методов кластеризации. Усвоив простые принципы, которые используются в этом алгоритме, вы получите знания, необходимые для понимания практически любого современного алгоритма кластеризации. Многие такие методы перечислены на следующем сайте, где представлен обзор CRAN-задач кластеризации: <http://cran.r-project.org/web/views/Cluster.html>.



По мере развития метода k -средних появилось множество реализаций этого алгоритма. Один из популярных подходов описан в статье: Hartigan J.A., Wong M.A. A k -means clustering algorithm. Applied Statistics, 1979. Vol. 28. P. 100–108.

Несмотря на то что с момента создания метода k -средних кластерные методы значительно усовершенствовались, это не означает, что метод k -средних устарел. Напротив, сегодня этот метод, возможно, более популярен, чем когда-либо. В табл. 9.1 перечислены причины, по которым метод k -средних все еще широко применяется.

Таблица 9.1

Преимущества	Недостатки
<p>Основан на простых принципах, которые можно объяснить без использования статистических терминов.</p> <p>Очень гибкий; достаточно простых настроек, чтобы устранить почти все его недостатки.</p> <p>Хорошо работает для решения многих практических задач</p>	<p>Не такой продвинутый, как более современные алгоритмы кластеризации.</p> <p>Из-за использования случайного выбора нет гарантии, что будет найден оптимальный набор кластеров.</p> <p>Требуется разумное предположение о том, сколько кластеров естественным образом существует в данных.</p> <p>Неидеален для несферических кластеров или кластеров с разной плотностью</p>

Если название метода k -средних кажется вам знакомым, вероятно, вы вспомнили алгоритм k ближайших соседей (k -NN), описанный в главе 3. Как вы скоро увидите, у метода k -средних и алгоритма k -NN гораздо больше общих черт, чем просто буква k в названии.

Алгоритм k -средних относит каждый из n примеров к одному из k кластеров, где k — число, которое определено заранее. Цель состоит в том, чтобы свести к минимуму различия между элементами каждого кластера и к максимуму — различия между кластерами.

Если k и n не очень малы, то вычислить оптимальные кластеры для всех возможных комбинаций примеров невозможно. В данном случае алгоритм использует эвристический процесс, который находит **локально оптимальные** решения. Это означает, что алгоритм начинает работу с некоторого начального приближенного разделения на кластеры и затем немного изменяет это разделение, чтобы увидеть, не улучшится ли однородность внутри кластеров в результате такого изменения.

Вскоре рассмотрим подробно этот процесс, но алгоритм, по существу, состоит из двух этапов. На первом этапе он присваивает каждый пример одному из k кластеров, образуя начальный набор. Затем алгоритм изменяет это разделение путем корректировки границ кластера в соответствии с примерами, которые в настоящее время попадают в кластер. Процесс изменения и разделения повторяется несколько раз, пока изменения не перестанут приводить к улучшениям разделения на кластеры. На этом этапе процесс останавливается и кластеры закрываются.



Благодаря эвристическому характеру метода k -средних можно в итоге получить несколько разных результатов за счет незначительного изменения начальных условий. Если результаты сильно различаются, это может указывать на проблему. Например, данные могут не иметь естественного разделения на группы, или значение k может быть выбрано неверно. С учетом этого имеет смысл выполнить кластерный анализ несколько раз, чтобы проверить надежность результатов. Чтобы увидеть, как на практике работает процесс разбиения на кластеры и внесения изменений, вернемся к примеру с конференцией на тему анализа данных. Несмотря на то что это простой пример, он наглядно демонстрирует основные принципы работы метода k -средних.

Использование расстояния для разбиения на кластеры и внесения изменений

Как и в случае с k -NN, метод k -средних обрабатывает значения признаков как координаты в многомерном пространстве. У данных о конференции существует только два признака, поэтому можно представить пространство объектов как двумерную диаграмму рассеяния, показанную ранее.

Алгоритм k -средних начинает работу с выбора в пространстве признаков k точек, которые будут служить центрами кластеров. Эти центры являются активаторами, которые будут побуждать остальные примеры встать на свои места. Часто точки выбираются путем выбора k случайных примеров из тренировочного набора данных. Поскольку мы собираемся идентифицировать три кластера, используя этот метод, то случайным образом выберем $k = 3$ точки. На рис. 9.3 эти точки обозначены звездочкой, треугольником и ромбом.

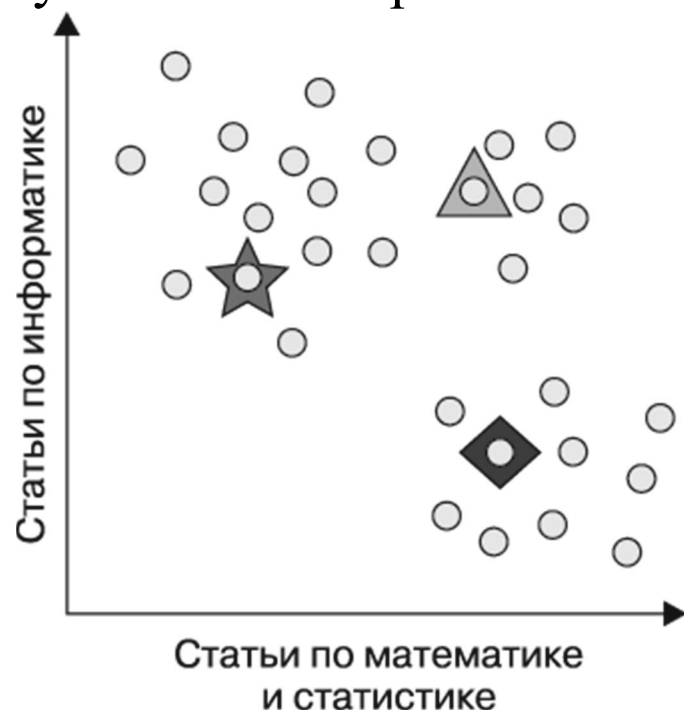


Рис. 9.3. Кластеризация методом k -средних начинается с выбора k случайных центров кластеров

Следует отметить, что на рис. 9.3 три центра кластеров расположены на большом расстоянии друг от друга, однако это не всегда так. Поскольку центры выбираются случайным образом, они могли бы с тем же успехом

оказаться по соседству. И поскольку алгоритм k -средних очень чувствителен к начальной позиции центров кластеров, это означает, что случайный выбор может существенно повлиять на итоговый набор кластеров.

Чтобы разрешить эту проблему, метод k -средних может быть дополнен различными методами выбора начальных центров. Например, в одном из вариантов выбираются случайные значения, встречающиеся в любом месте пространства признаков (а не только среди значений, присутствующих в данных). Другой вариант — вообще пропустить этот шаг — произвольно присвоить каждый пример какому-либо кластеру и сразу перейти к этапу внесения изменений. Каждый из этих подходов добавляет определенный сдвиг в окончательный набор кластеров, который можно будет использовать для улучшения результатов.



В 2007 году появился алгоритм под названием «метод k -средних++», в котором предлагается альтернативный способ выбора начальных центров кластеров. Предполагается, что это эффективный способ приблизиться к оптимальному кластерному решению, который заключается в уменьшении влияния случайного выбора. Подробнее об алгоритме k -средних++ читайте в статье: Arthur D., Vassilvitskii S. The advantages of careful seeding // Proceedings of the eighteenth annual ACM-SIAM symposium on discrete algorithms, 2007. P. 1027–1035.

После выбора начальных центров кластеров остальные примеры присваиваются тому кластеру, центр которого является ближайшим к данному примеру в соответствии с функцией расстояния. Как вы, вероятно, помните, мы рассмотрели функции расстояния при изучении алгоритма k -NN. Для метода k -средних обычно используется евклидово расстояние, но иногда вместо него используются манхэттенское расстояние или расстояние Минковского.

Напомним, что если n — число объектов, то формула евклидова расстояния между примером x и примером y выглядит следующим образом:

$$\text{dist}(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}.$$

Например, если сравнивать гостя, у которого пять публикаций по информатике и одна по математике, с гостем, у которого нет статей по информатике, но есть две статьи по математике, то можно вычислить расстояние в \mathbb{R} следующим образом:

```
> sqrt((5 - 0)^2 + (1 - 2)^2)[1] 5.09902
```

Используя эту функцию расстояния, можно вычислить расстояние между каждым примером и каждым центром кластера. Затем пример

присваивается тому кластеру, центр которого находится к нему ближе всего.



Следует помнить, что, поскольку мы вычисляем расстояния, все объекты должны быть числовыми, а значения — заранее нормализованными к стандартному диапазону. Для решения этой задачи полезно применить методы, представленные в главе 3.

Как показано на рис. 9.4, три центра кластеров разбивают примеры на три сегмента, обозначенные как *кластер А*, *кластер Б* и *кластер В*. Пунктирные линии показывают границы на *диаграмме Вороного*, образованной центрами кластеров. Диаграмма Вороного показывает области, расположенные ближе других к центру кластера; точка, в которой встречаются все три границы, является максимально удаленной от всех трех центров кластеров. Используя эти границы, легко увидеть области, созданные вокруг каждого из начальных центров формирования кластеров в методе *k*-средних.

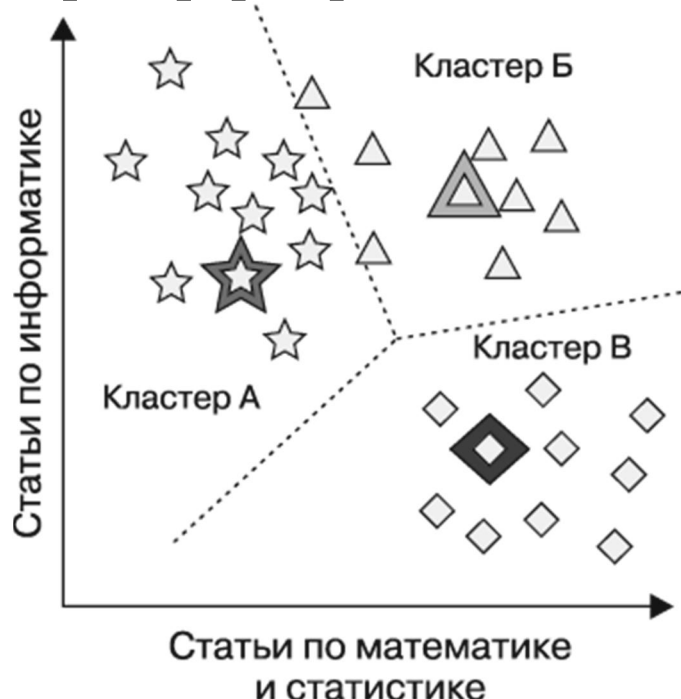


Рис. 9.4. Начальные центры кластеров образуют три группы ближайших точек

Теперь, когда начальный этап назначения завершен, алгоритм *k*-средних переходит к этапу внесения изменений. Первый шаг изменения кластеров — это перемещение начальных центров в новое положение, называемое *центроидом*. Центроид вычисляется как среднее положение точек, в настоящее время назначенных этому кластеру. На рис. 9.5 показано, что при перемещении центров кластеров к новым центроидам границы диаграммы Вороного также смещаются и точка, которая когда-то принадлежала кластеру Б (отмечена стрелкой), переносится в кластер А.

В результате этого переназначения алгоритм *k*-средних продолжит выполнение, и начнется следующий этап изменения. На рис. 9.6 показано, как выглядит картина после смещения центроидов кластера, изменения границ кластеров и переноса точек в новые кластеры (показано стрелками).

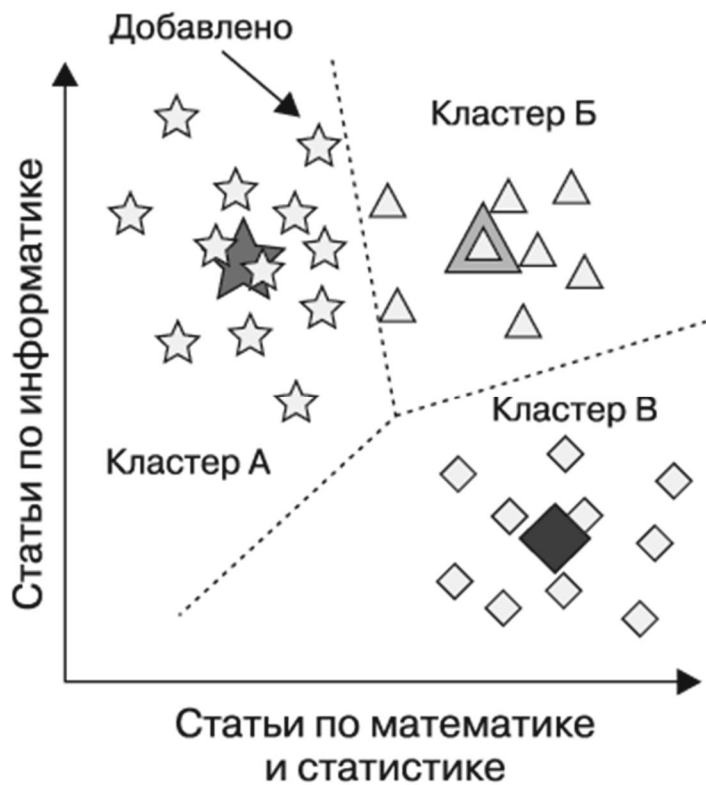


Рис. 9.5. На этапе внесения изменений смещаются центры кластеров, что приводит к переносу одной из точек в другой кластер

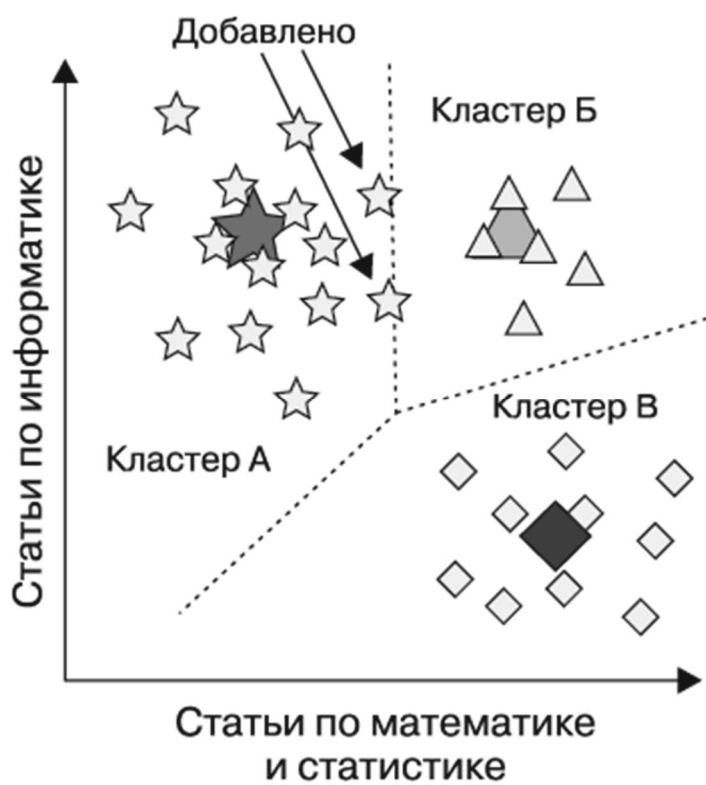


Рис. 9.6. После второго изменения еще две точки переносятся в кластер, центр которого расположен к ним ближе всего

Поскольку были перенесены еще две точки, должно произойти еще одно изменение с перемещением центроидов и обновлением границ кластеров. Однако, поскольку эти изменения не приводят к переносу точек, алгоритм *k*-средних останавливается. Разбиение по кластерам теперь является окончательным (рис. 9.7).



Рис. 9.7. Алгоритм кластеризации останавливается после внесения изменений, если это не привело к переносу точек в другие кластеры

Полученные в результате кластеры могут быть представлены одним из двух способов. Во-первых, можно просто сообщить о принадлежности каждого примера из набора данных кластеру А, Б или В. Во-вторых, можно сообщить координаты центроидов кластера, полученные после последнего изменения. Любой из этих вариантов представления отчета легко преобразовать в другой, вычисляя центроиды по координатам примеров из каждого кластера или же присваивая каждый пример тому кластеру, чей центроид находится ближе всего.

Выбор количества кластеров

Знакомясь с методом k -средних, мы узнали, что этот алгоритм чувствителен к случайно выбранным центрам кластеров. Действительно, если бы в предыдущем примере мы выбрали другие три начальные точки, то мы бы построили кластеры, которые разделяют данные не так, как ожидалось. Кроме того, метод k -средних чувствителен к количеству кластеров; выбор k требует поиска тонкого баланса. Слишком большое значение k повысит однородность кластеров, однако в то же время может привести к чрезмерной аппроксимации данных.

В идеале нужны *априорные* знания (предварительное предположение) о реально существующих группах, чтобы применить эту информацию для выбора количества кластеров. Например, при разбиении на группы фильмов можно начать с k , равного количеству жанров, которые рассматриваются в номинациях на премию «Оскар». В задаче о размещении участников конференции k может соответствовать количеству основных академических направлений исследований, представители которых приглашены на мероприятие.

Иногда количество кластеров продиктовано требованиями бизнеса или причинами анализа. Например, количество столов в зале заседаний

определяется исходя из того, сколько групп следует создать согласно списку участников конференции.

Вот еще один пример, где показано экономическое обоснование количества кластеров: если у отдела маркетинга хватает ресурсов только на три рекламные кампании, то, возможно, имеет смысл установить $k = 3$, чтобы каждый потенциальный клиент попал в одну из этих трех кампаний.

На случай, если никаких предварительных знаний нет, одно из эмпирических правил гласит: установить k равным квадратному корню из $(n / 2)$, где n — количество примеров в наборе данных. Однако в случае больших наборов данных это практическое правило может привести к неподъемному количеству кластеров. К счастью, существуют другие количественные методы, которые могут помочь в поиске подходящего количества кластеров для метода k -средних.

Технология «метод локтя» пытается определить, как изменяется однородность или неоднородность кластеров при различных значениях k . Как показано на рис. 9.8, однородность внутри кластеров, по идее, должна увеличиваться при добавлении новых кластеров; и наоборот, с увеличением количества кластеров неоднородность будет уменьшаться. Поскольку улучшения будут продолжаться до тех пор, пока каждый пример не окажется в своем отдельном кластере, цель состоит не в том, чтобы увеличивать однородность или уменьшать неоднородность до бесконечности, а скорее в том, чтобы найти k , при котором изменения становятся ниже некоторого порога. Такое значение k называется *точкой перегиба*, потому что диаграмма в этом месте имеет перегиб.

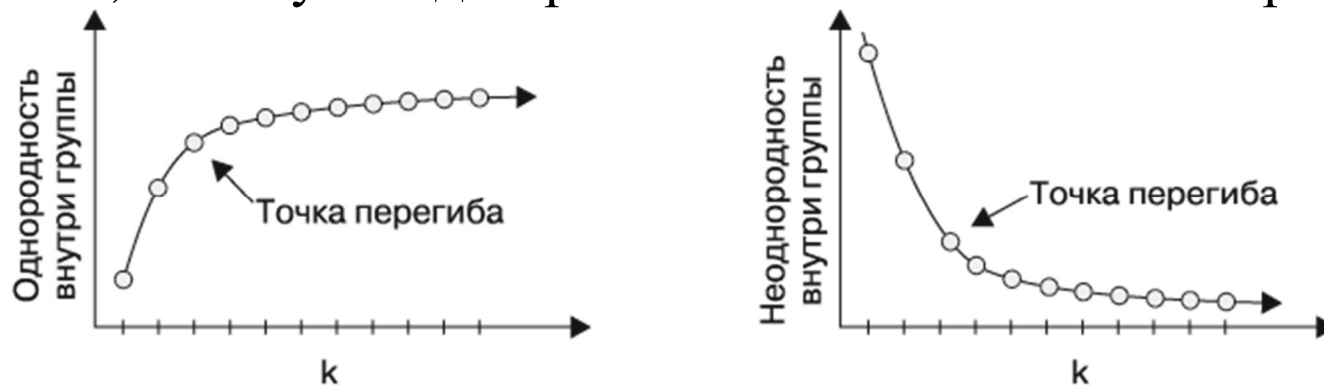


Рис. 9.8. Точка перегиба — это точка, в которой увеличение k приводит к относительно незначительным улучшениям

Существует множество статистических показателей, позволяющих измерить однородность и неоднородность внутри кластеров. Эти показатели можно использовать в методе локтя. Однако на практике итеративное тестирование большого количества значений k не всегда выполнимо. Отчасти так происходит потому, что кластеризация в случае больших наборов данных может занимать довольно много времени, а многократная кластеризация — еще больше. Кроме того, задачи, требующие точного оптимального набора кластеров, встречаются довольно редко. В большинстве задач кластеризации достаточно выбрать значение k для удобства, а не из-за строгих требований к эффективности.



Подробный обзор многочисленных и разнообразных показателей эффективности выбора кластеров представлен в статье: Halkidi M., Batistakis Y., Vazirgiannis M. On Clustering Validation Techniques // Journal of Intelligent Information Systems, 2001. Vol. 17. P. 107–145.

Уже сам процесс выбора k иногда приводит к интересным выводам. Наблюдая, как изменяются характеристики кластеров при изменении k , можно сделать вывод о том, где находятся естественные границы данных. Более плотные группы изменяются очень мало, в то время как менее однородные группы будут постоянно то распадаться, то вновь образовываться.

Как правило, иногда стоит потратить немного времени на то, чтобы получить правильное значение k . В следующем примере показано, как даже крупницы знаний предметной области, заимствованные из голливудского фильма, могут быть использованы для выбора такого k , при котором будут получены полезные и интересные кластеры. Поскольку кластеризация является неконтролируемой, на самом деле задача заключается в извлечении полезной информации — в знаниях, полученных из результатов работы алгоритма.

Сегментация рынка для подростков с использованием кластеризации методом k -средних

Общение с друзьями в *социальных сетях* (Social Networking Service, SNS), таких как Facebook, Tumblr и Instagram, стало для подростков всего мира обычным делом. Имея достаточное количество наличных денег, подростки являются желанной социально-демографической группой для компаний, которые продают закуски, напитки, электронику и средства гигиены.

Миллионы подростков, посещающих такие сайты, привлекли внимание маркетологов, стремящихся найти свою нишу на все более высококонкурентном рынке. Один из способов найти такую нишу — выявление среди подростков групп, имеющих схожие вкусы, чтобы клиенты, не заинтересованные в этих товарах, не получали рекламу, ориентированную на подростков. Например, скорее всего, будет трудно продать спортивную одежду тем подросткам, которые не интересуются спортом.

Исходя из информации на страницах подростков в социальных сетях, можно выделить группы с общими интересами, такими как спорт, религия или музыка. Кластеризация может автоматизировать процесс обнаружения естественных сегментов в этой социально-возрастной группе. Однако только нам решать, насколько эти кластеры интересны и как их можно использовать для рекламы. Пройдем этот процесс от начала до конца.

Шаг 1. Сбор данных

Для этого анализа мы будем использовать набор данных, представляющий собой случайную выборку, состоящую из данных о 30 000 американских школьниках, у которых в 2006 году были профили в популярной социальной сети. Чтобы сохранить анонимность пользователей, я не буду упоминать название этой сети. Скажу лишь, что на момент сбора данных эта сеть была популярным местом общения для подростков США. Поэтому разумно предположить, что их профили представляют собой довольно широкий срез данных об американских подростках на 2006 год.



Я собрал этот набор данных в рамках собственного социологического исследования подростковой идентичности для Университета Нотр-Дам. Если вы используете данные в исследовательских целях, пожалуйста, дайте ссылку на эту главу. Полный набор данных доступен в файле `snsdata.csv`. Для того чтобы выполнить пример из этой главы, не забудьте сохранить этот файл в своем рабочем каталоге R.

Данные выбирались равномерно по четырем годам окончания средней школы (с 2006 по 2009 год), так что в них представлены старшие и младшие классы, а также студенты первых и вторых курсов на момент сбора данных. С помощью автоматического веб-сканера были загружены полные тексты профилей из социальных сетей. Также были записаны пол, возраст и количество друзей в соцсети для каждого подростка.

Для разделения оставшегося содержимого страницы социальной сети на отдельные слова был использован инструмент текстового анализа. Из 500 слов, чаще всего появляющихся на всех страницах, были выбраны 36 слов для представления пяти категорий интересов, среди которых внешкольная деятельность, мода, религия, отношения и антиобщественное поведение. Эти 36 слов включают в себя и такие термины, как «футбол», «сексуальность», «поцелуи», «библия», «шопинг», «смерть» и «наркотики». Окончательный набор данных указывает на то, сколько раз каждое слово появляется в профиле социальной сети каждого человека.

Шаг 2. Исследование и подготовка данных

Для загрузки данных во фрейм данных можно использовать стандартные настройки `read.csv()`:

```
> teens <- read.csv("snsdata.csv")
```

Теперь кратко рассмотрим особенности данных. Первые несколько строк вывода функции `str()` выглядят следующим образом:

```
> str(teens)'data.frame':   30000 obs. of  40
variables:$ gradyear   : int 2006 2006 2006 2006
```

```

2006 2006 2006 2006 ...$ gender      : Factor w/ 2
levels "F", "M": 2 1 2 1 NA 1 1 2 ...$
age      : num 19 18.8 18.3 18.9 19 ...$
friends  : int 7 0 69 0 10 142 72 17 52 39 ...$
basketball : int 0 0 0 0 0 0 0 0 0 0 ...

```

Как мы и ожидали, данные охватывают 30 000 подростков и включают четыре переменные, соответствующие персональным характеристикам, и 36 слов, обозначающих интересы.

Заметили ли вы нечто странное, касающееся переменной пола (`gender`)? Если вы внимательно посмотрели, то, возможно, заметили значение `NA`, которое смотрится неуместно на фоне `1` и `2`. `NA` — это способ R сообщить нам о том, что в данной записи **отсутствует значение** — мы не знаем пол этого человека. До сих пор мы не имели дела с отсутствующими данными, но для многих видов анализа это может стать серьезной проблемой.

Посмотрим, насколько серьезна эта проблема в данном случае. Один из вариантов — использовать команду `table()` следующим образом:

```
> table(teens$gender)  F      M 22054  5222
```

Этот результат и говорит нам о количестве значений `F` (женщины) и `M` (мужчины), но функция `table()` исключает значения `NA`, вместо того чтобы рассматривать их как отдельную категорию. Чтобы включить в статистику значения `NA` (если они есть), нужно просто ввести дополнительный параметр:

```
> table(teens$gender, useNA =
"ifany")  F      M  <NA> 22054  5222  2724
```

Здесь мы видим, что 2724 записи (9 %) не имеют данных о гендере. Интересно, что, согласно этим данным, женщин в социальных сетях в четыре раза больше, чем мужчин. Это позволяет предположить, что мужчины менее, чем женщины, склонны использовать эту социальную сеть.

Если изучить другие переменные во фрейме данных, то обнаружится, что, кроме пола, отсутствующие значения есть только у переменной возраста (`age`). Для числовых данных количество отсутствующих значений можно узнать с помощью команды `summary()`:

```
> summary(teens$age)  Min.   1st
Qu.   Median   Mean   3rd
Qu.     Max.  NA's 3.086  16.310  17.290  17.990
      18.260  106.900  5086
```

В общей сложности в 5086 записях (17 %) возраст не указан. Также вызывает беспокойство тот факт, что минимальное и максимальное

значения представляются бессмысленными: маловероятно, чтобы трехлетний ребенок или 106-летний старик учились в средней школе. Чтобы эти запредельные значения не создавали проблем при анализе, нужно очистить данные, прежде чем двигаться дальше.

Разумный диапазон возрастов для старшеклассников — это те, кому уже исполнилось 13 и еще нет 20 лет. Любой возраст, выходящий за пределы этого диапазона, следует рассматривать как отсутствие данных — мы не можем доверять указанному возрасту. Чтобы перекодировать переменную `age`, можно воспользоваться функцией `ifelse()`, присваивая переменной `teen$age` значение `teen$age`, если возраст — не менее 13 и не выше 20 лет; в противном случае присваивается значение `NA`:

```
> teens$age <- ifelse(teens$age >= 13 &
teens$age < 20, teens$age,
NA)
```

Еще раз проверив результаты `summary()`, мы видим, что теперь возрастной диапазон соответствует распределению, которое намного больше похоже на реальную среднюю школу:

```
> summary(teens$age)  Min.   1st
Qu.   Median     Mean   3rd
Qu.    Max.   NA's13.03   16.30   17.26   17.25
18.22   20.00   5523
```

К сожалению, теперь у нас появилась еще большая проблема из-за отсутствующих данных. Нужно найти способ что-то сделать с этими значениями, прежде чем продолжить анализ.

Подготовка данных: фиктивные переменные вместо отсутствующих значений

Простое решение для обработки отсутствующих значений состоит в том, чтобы исключить все записи, в которых присутствуют не все значения. Однако, если представить себе последствия такого подхода, то вы дважды подумаете, прежде чем это сделать, — то, что это легко, еще не значит, что это хорошая идея! Проблема такого подхода состоит в том, что, даже если пропущено не так уж много значений, можно легко исключить большую часть данных.

Например, предположим, что в наших данных люди с гендерными значениями `NA` и те, кто не указал данные о возрасте, — совершенно разные люди. Это будет означать, что, исключая тех, у кого не указан пол или возраст, мы исключим $9 + 17 = 26$ % данных, то есть более 7500 записей. И это с учетом того, что данные отсутствуют только для двух переменных! Чем больше пропущенных значений в наборе данных, тем

выше вероятность того, что будут исключены все записи. Скоро у вас останется лишь небольшое подмножество данных или, еще хуже, оставшиеся записи будут систематически отличаться от удаленных, то есть не будут репрезентативными для всего населения.

Альтернативное решение для категориальных данных, таких как пол, состоит в том, чтобы рассматривать отсутствующее значение как отдельную категорию. Например, вместо того, чтобы ограничивать пол значениями «женский» и «мужской», мы можем ввести дополнительную категорию для неизвестного пола. Это позволит использовать фиктивное кодирование, описанное в главе 3.

Как вы помните, фиктивное кодирование означает создание специальной, фиктивной двоичной переменной, принимающей значения 1 или 0 для каждого уровня номинального признака, кроме одного, который используется в качестве контрольного. Причина, по которой одну из категорий можно исключить, состоит в том, что статус этой категории определяется по статусам остальных категорий. Например, если некто не является женщиной и его пол определен, то он должен быть мужчиной. Поэтому в данном случае нам нужно создать фиктивные переменные только для женского и неопределенного пола:

```
> teens$female <- ifelse(teens$gender == "F"
&                               !is.na(teens$gender),
1, 0)> teens$no_gender <-
ifelse(is.na(teens$gender), 1, 0)
```

Как следует из названия, функция `is.na()` проверяет, равен ли пол `NA`. Следовательно, первый оператор присваивает переменной `teens$female` значение 1, если пол равен `F` и не равен `NA`, в противном случае присваивается значение 0. Во втором операторе, если `is.na()` возвращает `TRUE`, то есть пол не определен, переменной `teens$no_gender` присваивается 1, в противном случае — 0. Чтобы убедиться, что мы все сделали правильно, сравним созданные нами фиктивные переменные с исходной переменной `gender`:

```
> table(teens$gender, useNA =
"ifany")      F      M      <NA> 22054      5222      2724>
table(teens$female, useNA =
"ifany")      0      17946      22054>
table(teens$no_gender, useNA =
"ifany")      0      127276      2724
```

Количество значений 1 для переменных `teens$female` и `teens$no_gender` соответствует

числу значений `F` и `NA` переменной `teens$gender` соответственно. Следовательно, нашей работе можно доверять.

Подготовка данных: подстановка отсутствующих значений
Теперь исключим 5523 отсутствующих значения возраста. Поскольку возраст является числом, то нет смысла создавать дополнительную категорию для неизвестных значений — как потом сравнивать «неизвестный» возраст с другими возрастами? Воспользуемся другой стратегией — *методом подстановок*. Этот метод подразумевает заполнение отсутствующих данных некими предположительно истинными значениями.

Можно ли придумать способ, как на основе данных из социальной сети сделать обоснованное предположение о возрасте подростка? Если вы подумали о годе окончания школы, вы на правильном пути. Большинство из тех, кто окончил школу в одном и том же году, родились в течение одного календарного года. Если определить типичный возраст для каждой группы выпускников, то получим достаточно разумную оценку возраста любого человека, который окончил школу в этом году.

Один из способов найти типичное значение — вычислить его среднее значение. Однако если применить к значениям возраста функцию `mean()`, как мы это делали в предыдущих задачах анализа, то возникнет проблема:

```
> mean(teens$age)[1] NA
```

Дело в том, что для вектора, содержащего недостающие данные, среднее значение не определено. Поскольку в наших данных о возрасте есть отсутствующие значения, `mean(teens$age)` возвращает `NA`. Это можно исправить, указав дополнительный параметр, позволяющий удалить пропущенные значения перед вычислением среднего:

```
> mean(teens$age, na.rm = TRUE)[1] 17.25243
```

Как видим, средний возраст студента в наших данных составляет примерно 17 лет. Но это только полдела: на самом деле нам нужен средний возраст для каждого года выпуска. Мы могли бы вычислить среднее значение четыре раза, но одним из преимуществ R является возможность избегать повторов. Здесь нам поможет функция `aggregate()`, которая вычисляет статистические показатели для подгрупп данных. В этом случае она вычисляет средний возраст по году выпуска после удаления значений `NA`:

```
> aggregate(data = teens, age ~ gradyear, mean,
na.rm =
TRUE)   gradyear      age1      2006      18.655862
```

```
2007 17.706173 2008 16.767704 200
9 15.81957
```

Как видим, средний возраст каждого года выпуска отличается примерно на один год. Это совершенно не удивительно, но является хорошим подтверждением того, что наши данные разумны.

Результатом `aggregate()` является фрейм данных. Поэтому, чтобы внести их в наши исходные данные, придется выполнить дополнительную работу. Или же можно использовать функцию `ave()`, которая возвращает вектор со средними значениями для группы, длина которого равна длине исходного вектора:

```
> ave_age <- ave(teens$age, teens$gradyear, FUN
= function(x) mean(x, na.rm =
TRUE))
```

Чтобы подставить эти средства вместо отсутствующих, нужно еще раз вызвать `ifelse()`, чтобы использовать значение `ave_age` только в тех случаях, когда исходное значение возраста равно `NA`:

```
teens$age <- ifelse(is.na(teens$age), ave_age,
teens$age)
```

Как видно по результатам `summary()`, пропущенных значений больше нет:

```
> summary(teens$age)  Min. 1st
Qu.  Median      Mean 3rd
Qu.    Max.13.03  16.28  17.24  17.24  18.2
1    20.00
```

Теперь, когда данные готовы для анализа, можно углубиться в самую интересную часть проекта. Посмотрим, окупались ли наши усилия.

Шаг 3. Обучение модели на данных

Чтобы разделить подростков по маркетинговым сегментам, воспользуемся реализацией алгоритма k-средних из пакета `stats`, который по умолчанию входит в состав установочного комплекта R. Если у вас нет этого пакета, вы можете его установить так же, как любой другой пакет, и загрузить с помощью команды `library(stats)`.

В различных R-пакетах есть множество функций k-средних, однако функция `kmeans()` из пакета `stats` остается весьма популярной и обеспечивает простую реализацию алгоритма.

Синтаксис кластеризации

Использование функции `kmeans()` из пакета `stats`

Поиск кластеров:

```
myclusters <- kmeans(mydata, k)
```

`mydata` – матрица или фрейм данных с примерами для кластеризации;

`k` – желаемое количество кластеров.

Функция возвращает объект кластеров, в котором хранится информация о созданных кластерах.

Исследование кластеров:

`myclusters$clusters` – вектор кластеров, сформированный функцией `kmeans()`;

`myclusters$centers` – матрица средних значений для каждого признака в кластерах;

`myclusters$size` – список количества примеров, попадающих в каждый кластер.

Пример:

```
teen_clusters <- kmeans(teens, 5)  
teens$cluster_id <- teen_clusters$cluster
```

Функция `kmeans()` принимает фрейм данных, содержащий только числовые данные, и параметр, указывающий желаемое количество кластеров. Если у вас есть и то и другое, то сам процесс построения модели прост. Проблема в том, что выбор правильного сочетания данных и числа кластеров — в некотором роде искусство; иногда это познается методом проб и ошибок.

Мы начнем кластерный анализ с рассмотрения всего 36 признаков, показывающих, сколько раз различные интересы появлялись в профилях социальной сети подростков. Для удобства построим фрейм данных, содержащий только эти признаки:

```
> interests <- teens[5:40]
```

Как было показано в главе 3, перед любым анализом, требующим вычисления расстояний, обычно выполняется нормализация или стандартизация по z-оценке таким образом, чтобы все признаки принадлежали одному и тому же диапазону. Это позволяет избежать проблемы, при которой одни признаки становятся важнее других только потому, что имеют более широкий диапазон значений.

При вычислении стандартной z-оценки признаки изменяются таким образом, что их среднее значение равно 0, а стандартное отклонение — 1. Интерпретация данных, измененная таким образом, может быть полезна в этом примере.

В частности, если некто в своем профиле трижды упоминает баскетбол, то, не имея дополнительной информации, нельзя сказать, означает ли это, что данный подросток любит баскетбол больше или меньше, чем его сверстники. Однако если z-оценка этого признака равна 3, значит, данный

пользователь соцсети упоминает баскетбол намного чаще, чем средний подросток.

Для того чтобы применить стандартизацию по z-оценке к фрейму данных `interests`, можно использовать функцию `scale()` в сочетании с `lapply()`. Поскольку `lapply()` возвращает матрицу, то эту матрицу необходимо преобразовать во фрейм данных с помощью функции `as.data.frame()` следующим образом:

```
> interests_z <-  
as.data.frame(lapply(interests, scale))
```

Чтобы убедиться, что преобразование выполнено правильно, сравним сводную статистику столбца `basketball` в старых и новых данных об интересах:

```
> summary(interests$basketball)      Min.   1st  
Qu.   Median     Mean  3rd  
Qu.   Max. 0.0000   0.0000   0.0000  0.2673   0.  
0000  24.0000>  
summary(interests_z$basketball)      Min.   1st  
Qu.   Median     Mean  3rd Qu.   Max. -0.3322  -  
0.3322  -0.3322   0.0000  -0.3322  29.4923
```

Как и ожидалось, в наборе данных `interests_z` признак `basketball` преобразован так, чтобы его среднее значение было равно 0, а диапазон охватывал значения больше и меньше 0. Теперь если этот признак меньше 0, то его можно интерпретировать так, что в соответствующем профиле соцсети баскетбол упоминается реже, чем в среднем для группы. Если значение больше 0, то этот человек упоминает баскетбол чаще среднего показателя.

Наше последнее решение включает в себя решение о том, сколько кластеров использовать для сегментации данных. Если кластеров окажется слишком много, то они могут оказаться чересчур специализированными и поэтому бесполезными; и наоборот, если кластеров слишком мало, это может привести к формированию гетерогенных групп. Не стесняйтесь экспериментировать с разными значениями k . Если вам не нравится результат, лучше выберите другое значение и начните все сначала.



Выбор количества кластеров упрощается, если вы знакомы с анализируемой демографической категорией. Имея представление о реальном количестве естественных групп, можно сэкономить время и избежать многих ошибок.

Чтобы проще было выбрать количество кластеров для этих данных, я предлагаю вам в помощь один из моих любимых фильмов — «Клуб

“Завтрак”» (The Breakfast Club) — подростковую комедию, снятую в 1985 году Джоном Хьюзом. В этом фильме пять героев-подростков — пять архетипов: «умник», «спортсмен», «безнадежный случай», «принцесса» и «преступник». Учитывая, что именно такие идентичности преобладают в популярном подростковом фильме, число 5 кажется разумной отправной точкой для k .

Чтобы использовать алгоритм k -средних для разделения данных об интересах подростков на пять кластеров, мы применим функцию `kmeans()` для фрейма данных `interests`. Поскольку в алгоритме k -средних используются случайные начальные точки, для того, чтобы ваши результаты соответствовали следующим примерам, используйте функцию `set.seed()`. Как вы, вероятно, помните из предыдущих глав, эта функция инициализирует генератор случайных чисел R определенной последовательностью. В отсутствие этой команды результаты могут меняться при каждом запуске алгоритма k -средних.

```
> set.seed(2345) > teen_clusters <-  
kmeans(interests_z, 5)
```

Результатом кластеризации методом k -средних является список `teen_clusters`, в котором хранятся свойства каждого из пяти кластеров. Изучим его внимательно и посмотрим, насколько хорошо алгоритм разделил данные об интересах подростков.



Если окажется, что ваши результаты отличаются от представленных здесь, проверьте, выполнили ли вы команду `set.seed(2345)` непосредственно перед `kmeans()`.

Шаг 4. Оценка эффективности модели

Оценка результатов кластеризации может быть несколько субъективной. В конце концов, успех или неудача модели зависит от того, насколько кластеры подходят для конкретной задачи. Поскольку этот анализ проводился, чтобы выявить группы подростков со схожими интересами в целях маркетинга, мы будем оценивать наш успех главным образом в качественном отношении. В других областях применения кластеризации могут потребоваться количественные показатели успешности.

Одним из основных способов оценки полезности множества кластеров является то, сколько примеров попало в каждую из групп. Если группы слишком большие или слишком маленькие, то едва ли от них будет много пользы. Чтобы вычислить размер кластеров, полученных с помощью `kmeans()`, используйте компонент `teen_clusters$size`:

```
>
```

```
teen_clusters$size[1]      871      600      5981      1034  
21514
```

Как видим, функция сформировала пять кластеров, как и было задано. В самый маленький из них попало 600 подростков (2 %), а в самый большой — 21,514 (72 %). Значительный разрыв между количеством подростков в самом большом и самом маленьком кластерах вызывает беспокойство, однако без более тщательного анализа этих групп мы не знаем, указывает ли это на проблему.

Может оказаться, что несоответствие размеров кластеров указывает на какие-то реальные обстоятельства — например, на большую группу подростков, имеющих общие интересы. Или же это может быть случайное совпадение, вызванное начальным выбором центров кластеров в методе k-средних. Чтобы в этом разобраться, нужно исследовать однородность каждого кластера.



Бывает так, что метод k-средних создает очень маленькие группы — иногда всего из одной точки. Такое может произойти, если один из начальных центров кластеризации оказался на значительном расстоянии от остальных данных. Не всегда ясно, следует ли рассматривать такие маленькие кластеры как истинный факт, представляющий группу крайних случаев, или же как проблему, вызванную случайным выбором. Если вы столкнетесь с этим явлением, то, возможно, стоит перезапустить алгоритм k-средних с другим случайным начальным числом и посмотреть, сохранится ли небольшой кластер при различных начальных точках. Для более глубокого исследования кластеризации мы рассмотрим координаты центроидов кластеров, используя компонент `teen_clusters$centers`, который для первых четырех интересов выглядит следующим образом:

```
>
```

```
teen_clusters$centers  basketball  football  
soccer  softball1  0.16001227  0.2364174  0.1  
0385512  0.072320212 -0.09195886  0.0652625 -  
0.09932124 -  
0.017394283  0.52755083  0.4873480  0.29778605  0  
.371788774  0.34081039  0.3593965  0.12722250  0.  
163846615 -0.16695523 -0.1641499 -0.09033520 -  
0.11367669
```

Строки выходных данных (помеченные цифрами от 1 до 5) соответствуют пяти кластерам, а числа в каждой строке означают среднее

значение в данном кластере для интересов, перечисленных в верхней части каждого столбца. Поскольку значения стандартизированы по z-оценке, то положительные числа означают превышение общего среднего уровня для всех подростков, а отрицательные — ниже среднего. Например, в третьей строке самое большое значение находится в столбце **basketball**, что означает, что в кластере **3** самый высокий из всех кластеров средний интерес к баскетболу.

Исследуя, находятся ли значения кластеров выше или ниже среднего уровня для каждой категории интересов, можно заметить закономерности, которые отличают кластеры друг от друга. На практике это означает вывод данных о центрах кластеров и поиск в этих данных каких-либо паттернов или экстремальных значений — почти как в головоломке по поиску слов, только в данном случае — чисел. На рис. 9.9 изображен снимок экрана с разноцветными выделениями паттернов, обнаруженных для каждого из пяти кластеров, для 19 из 36 интересов подростков.

```
> teen_clusters$centers
  basketball  football      soccer  softball  volleyball  swimming
1  0.16001227  0.2364174  0.10385512  0.07232021  0.18897158  0.23970234
2 -0.09195886  0.0652625 -0.09932124 -0.01739428 -0.06219308  0.03339844
3  0.52755083  0.4873480  0.29778605  0.37178877  0.37986175  0.29628671
4  0.34081039  0.3593965  0.12722250  0.16384661  0.11032200  0.26943332
5 -0.16695523 -0.1641499 -0.09033520 -0.11367669 -0.11682181 -0.10595448
cheerleading  baseball      tennis      sports      cute      sex
1  0.3931445  0.02993479  0.13532387  0.10257837  0.37884271  0.020042068
2 -0.1101103 -0.11487510  0.04062204 -0.09899231 -0.03265037 -0.042486141
3  0.3303485  0.35231971  0.14057808  0.32967130  0.54442929  0.002913623
4  0.1856664  0.27527088  0.10980958  0.79711920  0.47866008  2.028471066
5 -0.1136077 -0.10918483 -0.05097057 -0.13135334 -0.18878627 -0.097928345
  sexy      hot      kissed      dance      band      marching      music
1  0.11740551  0.41389104  0.06787768  0.22780899 -0.10257102 -0.10942590  0.1378306
2 -0.04329091 -0.03812345 -0.04554933  0.04573186  4.06726666  5.25757242  0.4981238
3  0.24040196  0.38551819 -0.03356121  0.45662534 -0.02120728 -0.10880541  0.2844999
4  0.51266080  0.31708549  2.97973077  0.45535061  0.38053621 -0.02014608  1.1367885
5 -0.09501817 -0.13810894 -0.13535855 -0.15932739 -0.12167214 -0.11098063 -0.1532006
```

Рис. 9.9. Чтобы различать кластеры, иногда полезно выделить паттерны

На основании этих данных об интересах уже можно определить некоторые особенности кластеров. Третий кластер отличается более высокими, чем средний уровень, интересами по всем видам спорта. В соответствии с архетипами из «Клуба “Завтрак”» это может быть группа «спортсменов». В первом кластере больше всего упоминаний слов «чирлидинг» и «горячий». Возможно, это те самые «принцессы»?

Продолжая исследовать кластеры подобным образом, можно составить таблицу, в которой будут перечислены доминирующие интересы каждой из групп. На рис. 9.10 представлены все кластеры в соответствии с особенностями, которые больше всего отличают каждый из них от остальных, и типом согласно фильму «Клуб “Завтрак”», который наиболее точно отражает характеристики данной группы.

Кластер 1 (N = 3376)	Кластер 2 (N = 601)	Кластер 3 (N = 1036)	Кластер 4 (N = 3279)	Кластер 5 (N = 21 708)
Плавание чирлидинг симпатичный сексуальный заводной танец платье волосы почта Холлистер Эберкромби шопинг одежда	Рок-группа поход музыка рок	Спорт секс сексуальный заводной целовал танец музыка рок-группа умирать смерть пьяный наркотики	Баскетбол футбол американский футбол волейбол бейсбол спорт Бог церковь Иисус Библия	???
Принцессы	Умники	Преступники	Спортсмены	Психи

Рис. 9.10. Таблица, которую можно использовать для перечисления важных свойств каждого кластера

Интересно, что пятый кластер выделяется своей неискл​ючительностью: интерес его членов ко всем выбранным видам деятельности ниже среднего. Эта группа также резко отличается от остальных большим количеством членов. Одно из возможных объяснений — эти пользователи соцсети создали профиль на сайте, но никогда не публиковали информации о своих интересах.



Предоставляя другим результаты анализа сегментации, иногда полезно использовать информационные метки, которые упрощают и отражают суть групп, — например, применяемую здесь типизацию по фильму «Клуб “Завтрак”». Риск создания таких меток заключается в том, что они могут скрывать особенности групп, создавая архетипы для их членов. Поскольку такие метки способны влиять на наше мышление, могут быть упущены важные закономерности, если воспринимать эти метки как исчерпывающую информацию.

Имея такую таблицу, менеджер по маркетингу получит четкое описание пяти типов подростков — посетителей сайта социальной сети.

Основываясь на их профилях, руководитель отдела маркетинга может продавать показы целевой рекламы компаниям, производящим товары, соответствующие одному или нескольким кластерам. В следующем разделе вы увидите, как в этих целях можно применять метки кластера к исходной социально-демографической группе.

Шаг 5. Повышение эффективности модели

Поскольку при кластеризации создается новая информация, эффективность алгоритма кластеризации зависит как минимум от качества самих кластеров и от того, что алгоритм делает с этой информацией. Из

предыдущего раздела видно, что пять кластеров дали полезную и новую информацию об интересах подростков. В этом отношении алгоритм, похоже, работает довольно хорошо. Поэтому теперь нужно направить усилия на то, чтобы претворить эти идеи в жизнь.

Для начала нужно применить кластеры к полному набору данных. Объект `teen_clusters`, созданный функцией `kmeans()`, включает в себя компонент `cluster`, в котором все 30 000 человек из образцового набора данных разбиты по кластерам. Мы можем добавить эту информацию как новый столбец во фрейме данных `teens` с помощью следующей команды:

```
> teens$cluster <- teen_clusters$cluster
```

Получив эти новые данные, можно проверить, как назначение кластера связано с отдельными характеристиками. Например, вот персональные данные для первых пяти подростков из данных о социальной сети:

```
> teens[1:5, c("cluster", "gender", "age",  
"friends")]  
  cluster gender    age friends1  
      5      M 18.982    72      3      F  
18.801    03      5      M 18.335    69  
4      5      F 18.875    05      4    <N  
A> 18.995    10
```

Используя функцию `aggregate()`, можно также узнать демографические характеристики кластеров. Как видим, средний возраст не сильно изменяется от кластера к кластеру, что не так уж удивительно, ведь все это подростки, которые еще не перешли в старшие классы. Картина выглядит следующим образом:

```
> aggregate(data = teens, age ~ cluster,  
mean)  cluster    age1    1 16.864972  
      2 17.390373    3 17.076564    4 17.119  
575      5 17.29849
```

Зато есть существенные различия в распределении женщин по кластерам — очень интересная находка, поскольку при формировании кластеров мы не использовали гендерные данные, однако по принадлежности к кластеру можно прогнозировать пол:

```
> aggregate(data = teens, female ~ cluster,  
mean)  cluster    female1    1 0.83811712  
      2 0.72500003    3 0.83781984    4 0.  
80270795      5 0.6994515
```

Напомню, что в целом около 74 % пользователей SNS — женщины. Первая группа, так называемые «принцессы», насчитывает почти 84 % женщин, в то время как вторая и пятая — всего около 70 %. Эти различия

означают, что существует разница в интересах, которые мальчики и девочки-подростки обсуждают на своих страницах в социальных сетях.

С учетом успешного прогнозирования пола, по принадлежности к кластерам, вероятно, можно также спрогнозировать количество друзей у пользователей. Похоже, эта гипотеза подтверждается следующими данными:

```
> aggregate(data = teens, friends ~ cluster,
mean)
 cluster  friends1      1  41.430542
      2  32.573333      3  37.161854      4  30.502
905      5  27.70052
```

В среднем больше всего друзей у «принцесс» (41,4); за ними следуют «спортсмены» (37,2) и «умники» (32,6). На нижнем конце шкалы находятся «преступники» (30,5) и «безнадежные случаи» (27,7). Как и в случае с полом, зависимость количества друзей у подростка от принадлежности к кластеру весьма примечательна, поскольку мы не использовали данные о дружбе в качестве исходных данных алгоритма кластеризации. Интересно также, что количество друзей, по-видимому, связано со стереотипом о популярности разных групп в старших классах: у популярных групп, как правило, больше друзей.

Взаимосвязь между членством в группе, полом и количеством друзей предполагает, что кластеры могут быть полезны для прогнозирования поведения. Проверка возможностей кластеров для такого прогнозирования позволит упростить работу с ними после передачи в отдел маркетинга, что в итоге повысит эффективность алгоритма.

Резюме

Сделанные выводы подтверждают популярную поговорку «Рыбак рыбака видит издалека». Используя методы машинного обучения для объединения подростков в группы со схожими интересами, мы разработали типологию подростковой идентичности, которая определяет такие характеристики, как пол и количество друзей. Эти же методы можно применять в других областях, получая схожие результаты.

В данной главе охвачены лишь основные принципы кластеризации. Существует много вариантов алгоритма k-средних, а также много других алгоритмов кластеризации, которые приносят в задачу свои сдвиги и эвристики. Основные принципы, с которыми вы познакомились в этой главе, позволят вам освоить названные методы кластеризации и применять их для решения новых задач.

В следующей главе мы начнем изучать методы измерения эффективности алгоритмов обучения, применимые ко многим задачам машинного обучения. Мы все время уделяли внимание оценке эффективности обучения алгоритмов, однако для достижения

максимальной степени эффективности очень важно уметь определять и измерять ее в самых строгих условиях.

10. Оценка эффективности модели

В те времена, когда лишь богатые могли позволить себе образование, знания студентов не оценивали с помощью тестов и экзаменов.

Тестировали учителей, чтобы родители могли знать, достаточно ли хорошо те учат их детей, оправданно ли жалование преподавателей. Сегодня, разумеется, все по-другому. Теперь оценки нужны для того, чтобы отличать тех, кто учится хорошо, от тех, кто учится плохо, выделять перспективных студентов и т.д.

Учитывая важность этого процесса, для разработки точной системы оценок учащихся приложены большие усилия. Справедливые системы оценивания состоят из огромного количества вопросов, которые охватывают широкий круг тем и позволяют отличить истинные знания от случайных ответов. Хорошая система оценивания также заставляет студентов решать задачи, с которыми они прежде никогда не сталкивались. Таким образом, правильные ответы позволяют дать общую оценку знаниям.

Процесс оценки алгоритмов машинного обучения очень похож на процесс оценки студентов. Поскольку у каждого алгоритма есть свои сильные и слабые стороны, тесты должны отличать один алгоритм от другого. Также важно понимать, как алгоритмы обучения будут работать с новыми данными.

В этой главе представлена следующая информация, необходимая для оценки эффективности машинного обучения:

- причины, по которым одной лишь точности прогноза может оказаться недостаточно для измерения эффективности, а также показатели эффективности, которые можно было бы использовать;
- методы, гарантирующие, что показатели эффективности разумно отражают способность модели прогнозировать новые варианты;
- способы реализации этих полезных показателей и методов в применении к моделям прогнозирования, описанным в предыдущих главах, в среде R.

Подобно тому как наилучший способ изучить некую тему — это начать ее преподавать кому-то другому, процесс обучения и оценки эффективности машинного обучения даст вам более глубокое понимание уже изученных нами методов.

Измерение эффективности классификации

В предыдущих главах мы измеряли точность классификатора как число правильных прогнозов, разделенное на общее количество прогнозов. Таким образом мы получали долю случаев, в которых алгоритм обучения

оказывался прав. Предположим, у нас есть классификатор новорожденных, определяющий, являются ли они носителями излечимого, но потенциально фатального генетического дефекта. Классификатор дает верный прогноз для 99 990 из 100 000 младенцев. Это будет означать 99,99%-ную точность и частоту ошибок всего 0,01 %.

На первый взгляд, это очень ценный классификатор. Однако было бы разумно собрать дополнительную информацию, прежде чем доверить такому тесту жизнь ребенка. Что, если генетический дефект обнаруживается только у 10 из каждых 100 000 детей? Тест, который неизменно прогнозирует отсутствие дефекта, будет правильным для 99,99 % всех случаев, но неверным для 100 % самых важных случаев. Другими словами, даже если классификатор чрезвычайно точен, он будет не особенно полезен для предотвращения излечимых врожденных дефектов.



Это одно из последствий проблемы дисбаланса классов, возникающей, когда подавляющее большинство записей в наборе данных относятся к одному классу.

Есть много способов измерить эффективность классификатора, однако наилучший показатель — тот, который определяет, является ли классификатор успешным, когда применяется по прямому назначению. Очень важно определить показатели эффективности, которые будут показывать полезность, а не просто точность. Для этого мы рассмотрим различные альтернативные показатели эффективности, полученные из матрицы несоответствий. Однако сначала рассмотрим, как подготовить классификатор для оценки.

Прогнозы классификатора

Цель оценки эффективности модели классификации состоит в том, чтобы понять, как ее эффективность будет экстраполироваться на будущие случаи. Поскольку, как правило, невозможно протестировать неизвестную модель в реальной среде, обычно для моделирования будущих условий модель выполняет классификацию примеров из набора данных, похожих на те, которые предполагается обрабатывать в будущем. Анализируя ответы алгоритма обучения, можно узнать о его достоинствах и недостатках.

В предыдущих главах мы уже оценивали эффективность классификаторов, однако стоит подумать о том, какие данные для этого есть в нашем распоряжении:

- фактические значения класса;
- прогнозируемые значения класса;
- расчетная вероятность правильного прогноза.

Фактические и прогнозируемые значения класса могут быть очевидными, однако именно они являются ключевыми в оценке эффективности. Подобно тому как учитель при проверке ответов ученика использует список правильных ответов, нам тоже нужно знать правильный ответ при оценке прогнозов алгоритма обучения. Цель заключается в том, чтобы сформировать два вектора данных: один с правильными (фактическими) значениями класса, а второй — с прогнозируемыми. Оба вектора должны иметь одинаковое количество значений, хранящихся в одинаковом порядке. Прогнозируемые и фактические значения могут быть сохранены в виде отдельных R-векторов или как столбцы одного R-фрейма данных.

Получить эти данные легко. Фактические значения класса поступают непосредственно из целевого столбца тестового набора данных. Прогнозируемые значения класса получаются от классификатора, который построен на тренировочных данных и затем применяется к тестовым данным. Для большинства пакетов машинного обучения это включает в себя применение функции `predict()` к объекту модели и фрейму тестовых данных, например: `predictions <- predict(model, test_data)`.

До сих пор мы рассматривали только прогнозы классификации с использованием этих двух векторов данных, но большинство моделей предоставляют и другую полезную информацию. Даже если классификатор делает только один прогноз для каждого примера, в отношении одних решений он может быть более уверенным, чем в отношении других. Например, классификатор может быть на 99 % уверен, что СМС со словами «бесплатно» и «рингтоны» является спамом, но только на 51 % уверен, что СМС со словами «сегодня вечером» является спамом. В обоих случаях классификатор определяет сообщение как спам, но гораздо увереннее в одном решении, чем в другом (рис. 10.1).

Изучение этих вероятностей прогнозов, вычисленных самим алгоритмом, дает полезную информацию для оценки эффективности модели. Если две модели допускают одинаковое количество ошибок, но одна из них способна более точно оценить степень своей неопределенности, то такая модель более разумна. В идеале было бы хорошо построить такой алгоритм обучения, который был бы очень уверен в себе, делая правильный прогноз, и не уверен в случае сомнений. Ключевой частью оценки эффективности модели является баланс между уверенностью и осторожностью.



Рис. 10.1. Уверенность алгоритмов обучения в своих прогнозах может быть разной, даже если они обучаются на одних и тех же данных

Вызов функции для вычисления внутренней вероятности прогноза зависит от конкретного R-пакета. Как правило, для большинства классификаторов функция `predict()` позволяет использовать дополнительный параметр, чтобы указать желаемый тип прогнозирования. Для того чтобы получить один прогнозируемый класс, такой как `spam` или `ham`, обычно выбирается параметр `type="class"`. Чтобы получить вероятность прогнозирования, необходимо присвоить параметру `type` значение `"prob"`, `"posterior"`, `"raw"` или `"probability"`, в зависимости от используемого классификатора.



Все классификаторы, описанные в этой книге, позволяют вычислять вероятность прогноза. Параметр `type` включен в описание синтаксиса каждой из моделей.

Например, для того, чтобы получить прогноз вероятности для классификатора `C5.0`, о котором шла речь в главе 5, нужно использовать функцию `predict()` с параметром `type="prob"` следующим образом:

```
> predicted_prob <- predict(credit_model,
credit_test, type = "prob")
```

Для того чтобы вычислить вероятности прогнозов наивного байесовского алгоритма для модели классификации спама в СМС, построенной в главе 4, нужно использовать функцию `predict()` с параметром `type="raw"`:

```
> sms_test_prob <- predict(sms_classifier,
sms_test, type = "raw")
```

В большинстве случаев функция `predict()` возвращает вероятность для каждой категории результата. Например, в случае модели с двумя

вариантами, такой как классификатор СМС, прогнозируемые вероятности могут быть сохранены в виде матрицы или фрейма данных:

```
>
head(sms_test_prob)
      ham      sp
am[1,] 9.999995e-01 4.565938e-
07[2,] 9.999995e-01 4.540489e-
07[3,] 9.998418e-01 1.582360e-
04[4,] 9.999578e-01 4.223125e-
05[5,] 4.816137e-
10 1.000000e+00[6,] 9.997970e-01 2.030033e-04
```

Каждая строка здесь показывает прогнозируемую вероятность классификации писем как спам и не спам. Согласно правилам вычисления вероятности сумма значений в каждой строке должна быть равна 1, потому что это взаимоисключающие и исчерпывающие результаты. Учитывая эти данные, при построении набора данных для оценки эффективности модели необходимо убедиться, что выбрана только вероятность для интересующего вас класса. Для удобства в процессе оценки иногда полезно создать фрейм данных, содержащий прогнозируемый класс, фактический класс и прогнозируемую вероятность интересующего класса.



Для краткости изложения были опущены операции, необходимые для построения набора данных для оценки эффективности, но они включены в код этой главы. Чтобы выполнить приведенные здесь примеры, скачайте файл `sms_results.csv` и поместите его во фрейм данных с помощью команды `sms_results <- read.csv("sms_results.csv")`.

Структура фрейма данных `sms_results` очень проста. Он содержит четыре вектора из 1390 значений: один вектор с реальными типами СМС (спам или не спам), второй — с прогнозируемыми типами, полученными с помощью наивного байесовского алгоритма, а третий и четвертый векторы содержат вероятности того, что данное сообщение является спамом или не спамом соответственно:

```
>
head(sms_results)
  actual_type predict_type prob
b_spam      1.000002      ham      ham      0.000000
00000      1.000003      ham      ham      0.000000
016      0.999844      ham      ham      0.000004
04      0.999965      spam      spam      1.000000
0      0.000006      ham      ham      0.000200
      0.99980
```

Для этих шести тестовых примеров спрогнозированные и реальные типы СМС совпадают; модель правильно спрогнозировала их статус. Кроме того, вероятности предсказания говорят о том, что модель была чрезвычайно уверена в своих прогнозах, поскольку все вероятности близки к 0 или 1.

Что происходит, если вероятность прогнозируемых значений отстоит дальше от 0 или 1? Используя функцию `subset()`, можно найти несколько таких записей. В следующем выводе этой функции показаны тестовые случаи, в которых модель оценивала вероятность спама от 40 до 60 %:

```
> head(subset(sms_results, prob_spam > 0.40 &
prob_spam <
0.60))
```

		actual_type	predict_type	prob_spam
prob_ham377		spam	ham	0.47536
	0.52464717	ham	spam	0.561
88	0.438121311	ham	spam	0.5
7917	0.42083			

По собственной оценке модели, в этих случаях вероятность правильного прогноза практически не отличается от того, что мы могли бы получить, подбрасывая монетку. Тем не менее все три прогноза оказались неверными — неудачный результат. Рассмотрим еще несколько случаев, когда алгоритм ошибся:

```
> head(subset(sms_results, actual_type !=
predict_type))
```



		actual_type	predict_type	prob
_spam	prob_ham53	spam	ham	0
.00071	0.9992959	spam	ham	
0.00156	0.9984473	spam	ham	
0.01708	0.9829276	spam	ham	
0.00851	0.99149184	spam	ham	
0.01243	0.98757332	spam	ham	
0.00003	0.99997			

Эти случаи иллюстрируют важный факт: модель может быть чрезвычайно уверенной в своих прогнозах и при этом давать совершенно ошибочные результаты. Все шесть тестовых случаев были спамом, однако, по мнению классификатора, имели не менее 98%-ную вероятность оказаться обычными письмами.

С учетом таких ошибок можно ли по-прежнему считать такую модель полезной? Чтобы ответить на вопрос, нужно применить к этим данным оценки эффективности различные способы измерения ошибок. В сущности, многие такие измерения выполняются с помощью инструмента, который мы уже активно использовали в предыдущих главах.

Анализ матриц несоответствий

Матрица несоответствий — это таблица, которая классифицирует прогнозы в зависимости от того, соответствуют ли они фактическим значениям. Одно из измерений таблицы — это возможные категории прогнозируемых значений, а второе — те же категории для реальных значений. До сих пор нам встречались только матрицы несоответствий размерностью 2×2 , однако можно построить матрицу для моделей, которые прогнозируют любое количество значений класса. На рис. 10.2 изображена уже знакомая нам матрица несоответствий для двухклассовой двоичной модели, а также матрица несоответствий размерностью 3×3 для трехклассовой модели.

Если прогнозируемое значение совпадает с реальным, то классификация правильная. Правильные прогнозы размещаются в матрице несоответствий по диагонали (обозначаются как ). Остальные ячейки матрицы (обозначаются ) соответствуют случаям, когда прогнозируемое значение отличается от реального. Это неверные прогнозы. Показатели эффективности для моделей классификации основаны на количестве прогнозов, попадающих на диагональ и за ее пределы (рис. 10.2).

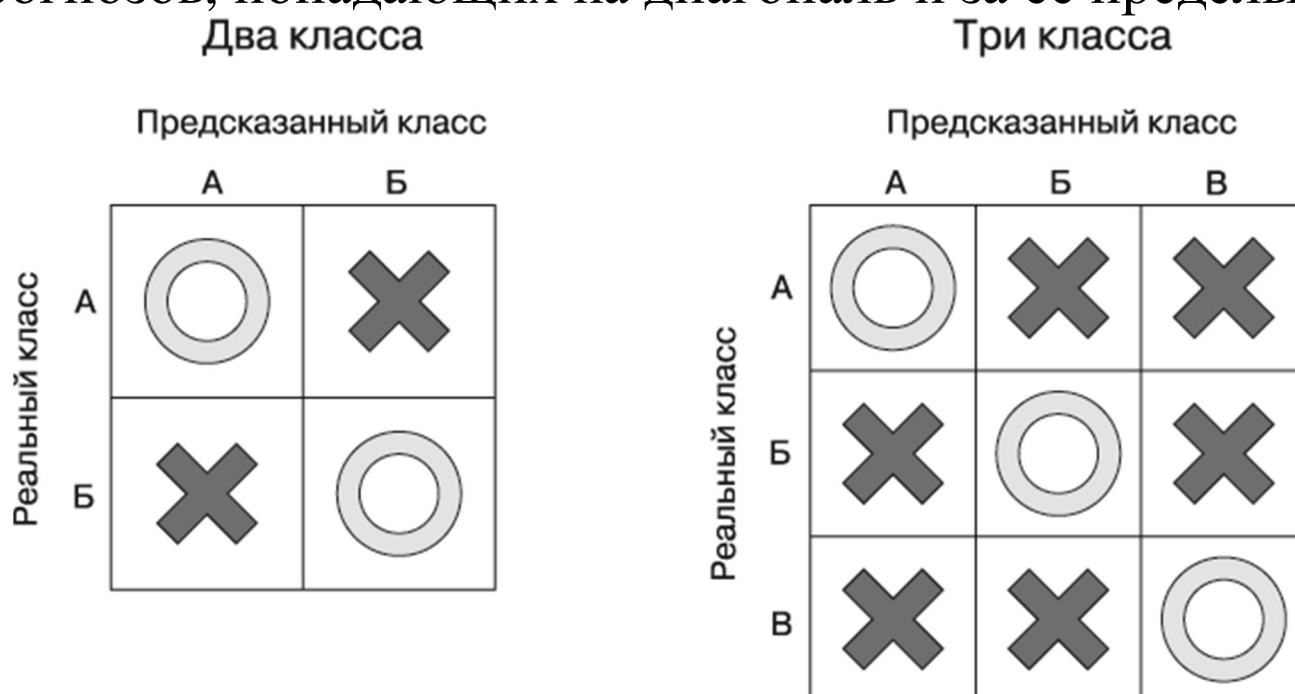


Рис. 10.2. Матрицы несоответствий позволяют подсчитать случаи, когда прогнозируемый класс совпадает или не совпадает с реальным

Наиболее распространенные показатели эффективности учитывают способность модели отличать один класс от остальных. При этом интересующий класс называется *положительным*, а остальные — *отрицательными*.



Термины «положительный» и «отрицательный» не носят характер оценочного суждения (то есть не означают «хороший» или «плохой») и не обязательно означают наличие или отсутствие результата (например, врожденный дефект или его отсутствие). Выбор положительного результата может быть даже произвольным, как в тех случаях, когда

модель прогнозирует такие категории, как солнечная или дождливая погода или собака либо кошка.

Зависимость между прогнозами положительного и отрицательного класса может быть представлена в виде матрицы несоответствий 2×2 , в которой показано, относятся ли прогнозы к одной из четырех категорий:

- *истинно положительная* (True Positive, TP) — правильно классифицируется как интересующий класс;
- *истинно отрицательная* (True Negative, TN) — правильно классифицирован как не относящийся к интересующему классу;
- *ложноположительная* (False Positive, FP) — неправильно классифицируется как интересующий класс;
- *ложноотрицательная* (False Negative, FN) — неправильно классифицирован как не относящийся к интересующему классу.

Для классификатора спама положительным классом является **spam**, так как именно этот результат мы стремимся обнаружить. Затем можно составить матрицу несоответствий (рис. 10.3).

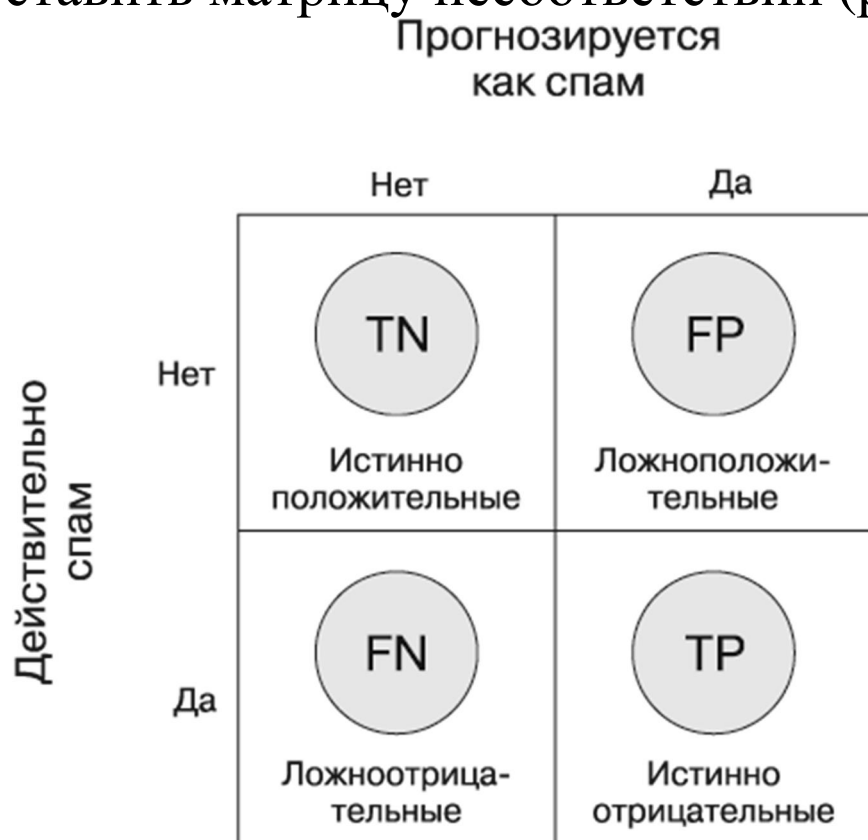


Рис. 10.3. Различие между положительными и отрицательными классами добавляет новые подробности к матрице несоответствий

Подобная матрица несоответствий является основой для многих важнейших показателей эффективности модели. В следующем разделе используем эту матрицу, чтобы лучше понять, что именно подразумевается под точностью.

Использование матриц несоответствий для измерения эффективности

С помощью матрицы несоответствий 2×2 можно формализовать определение точности прогнозирования (accuracy) (иногда называемое *коэффициентом успешности*):

$$\text{accuracy} = \frac{TP+TN}{TP+TN+FP+FN}.$$

В этой формуле TP, TN, FP и FN обозначают количество случаев, когда прогнозы модели попадают в каждую из соответствующих категорий. Таким образом, точность представляет собой отношение суммы истинно положительных и истинно отрицательных значений к общему количеству прогнозов.

Коэффициент ошибок (error rate), или доля неправильно классифицированных примеров, определяется так:

$$\text{error rate} = \frac{FP+FN}{TP+TN+FP+FN} = 1 - \text{accuracy}.$$

Обратите внимание, что коэффициент ошибок можно вычислить как 1 минус точность. Интуитивно это понятно: если модель является верной в 95 % случаев, то в 5 % случаев она неверна.

Простой способ преобразовать таблицу прогнозов классификатора в матрицу несоответствий — использовать R-функцию `table()`. Команда для создания матрицы несоответствий для данных об СМС показана ниже. Значения из этой таблицы впоследствии можно использовать для определения точности и других статистических показателей:

```
> table(sms_results$actual_type,
sms_results$predict_type)      ham  spam  ham
1203      4  spam    31  152
```

Если вы хотите создать матрицу путаницы с более информативными результатами, то можно воспользоваться функцией `CrossTable()` из пакета `gmodels`, которая предлагает более настраиваемое решение. Возможно, помните, впервые мы использовали эту функцию в главе 2. Если вы тогда не установили этот пакет, то сейчас вам нужно будет сделать это с помощью команды `install.packages("gmodels")`.

По умолчанию функция `CrossTable()` выводит в каждой ячейке пропорции, которые указывают на процентное соотношение относительно количества ячеек в строках, столбцах и в таблице в целом. В выводе функции также приводится общее количество строк и столбцов. Как показано в следующем коде, по синтаксису

функция `CrossTable()` похожа на `table()`:

```
> library(gmodels)>
CrossTable(sms_results$actual_type,
sms_results$predict_type)
```

В результате получается матрица несоответствий с множеством дополнительных деталей:

Cell Contents	
	N
Chi-square contribution	
N / Row Total	
N / Col Total	
N / Table Total	

Total Observations in Table: 1390

sms_results\$actual_type	sms_results\$predict_type		Row Total
	ham	spam	
ham	1203 16.128 0.997 0.975 0.865	4 127.580 0.003 0.026 0.003	1207 0.868
spam	31 106.377 0.169 0.025 0.022	152 841.470 0.831 0.974 0.109	183 0.132
Column Total	1234 0.888	156 0.112	1390

Мы уже использовали функцию `CrossTable()` в нескольких предыдущих главах, так что вы должны быть знакомы с результатами ее работы. Если вы забудете, как их интерпретировать, просто обратитесь к ее легенде (помеченной как `CellContents`), в которой дается определение всех чисел в ячейках таблицы.

Мы можем использовать матрицу несоответствий для определения точности и частоты ошибок. Поскольку точность равна $(TP + TN) / (TP + TN + FP + FN)$, то ее можно вычислить следующим образом:

```
> (152 + 1203) / (152 + 1203 + 4 + 31)[1]
0.9748201
```

Мы также можем вычислить коэффициент ошибок как $(FP + FN) / (TP + TN + FP + FN)$:

```
> (4 + 31) / (152 + 1203 + 4 + 31)[1]
0.02517986
```

Или как 1 минус точность:

```
> 1 - 0.9748201[1] 0.0251799
```

Эти вычисления могут показаться простыми, однако важно понять, каким образом компоненты матрицы несоответствий связаны между собой. В следующем разделе вы увидите, как можно по-разному комбинировать одни и те же элементы этой матрицы, чтобы получать различные дополнительные показатели эффективности.

Не только точность: другие показатели эффективности

Существует бесчисленное множество показателей эффективности, разработанных для конкретных целей и используемых в самых разных областях: в медицине, для поиска информации, в маркетинге, теории обнаружения сигналов и т.п. Их исчерпывающее описание заняло бы сотни страниц, что в данной книге невозможно. Рассмотрим лишь некоторые из самых полезных показателей, наиболее часто упоминаемых в литературе по машинному обучению.

Функции для вычисления многих таких показателей эффективности вы найдете в пакете классификации и регрессионного обучения `caret` от Макса Куна (Max Kuhn). В этот пакет входит большое количество инструментов для подготовки, обучения, оценки и визуализации моделей машинного обучения и данных. Мы используем этот пакет не только здесь, но и — более активно — в главе 11. Прежде чем продолжить чтение, вам потребуется установить этот пакет с помощью команды `install.packages("caret")`.



Подробнее о пакете `caret` читайте в статье: Kuhn M. Building Predictive Models in R Using the `caret` Package // Journal of Statistical Software, 2008. Vol. 28.

В пакет `caret` входит еще одна функция для построения матрицы несоответствий. Как показано в следующей команде, ее синтаксис похож на `table()`, но с небольшим отличием. Поскольку `caret` вычисляет показатели эффективности модели, отражающие ее способность классифицировать положительный класс, следует указывать параметр `positive`. Поэтому, поскольку классификатор СМС предназначен для обнаружения спама, в данном примере будем присваивать этому параметру значение `positive="spam"`:

```
> library(caret)>
confusionMatrix(sms_results$predict_type, sms_
results$actual_type, positive = "spam")
```

В результате получим следующий вывод:

Confusion Matrix and Statistics

```
Reference
Prediction ham spam
ham 1203 31
spam 4 152
```

```
Accuracy : 0.9748
95% CI : (0.9652, 0.9824)
No Information Rate : 0.8683
P-Value [Acc > NIR] : < 2.2e-16
```

```
Kappa : 0.8825
McNemar's Test P-Value : 1.109e-05
```

```
Sensitivity : 0.8306
Specificity : 0.9967
Pos Pred Value : 0.9744
Neg Pred Value : 0.9749
Prevalence : 0.1317
Detection Rate : 0.1094
Detection Prevalence : 0.1122
Balanced Accuracy : 0.9136
```

```
'Positive' Class : spam
```

В верхней части этих результатов размещается матрица несоответствий, очень похожая на матрицу, построенную функцией `table()`, только транспонированную. Здесь также есть набор показателей эффективности. Некоторые из них, такие как точность, вам уже знакомы, а другие встречаются впервые. Рассмотрим наиболее важные из показателей.

Каппа-статистика

Каппа-статистика (обозначенная выше в результатах функции `confusionMatrix` как `Kappa`) корректирует значение точности, учитывая то, что правильный прогноз может быть выдан случайно. Это особенно важно для наборов данных с серьезным дисбалансом классов, поскольку классификатор может получить высокое значение точности только потому, что случайно угадывает наиболее частый класс. Высокое значение каппа-статистики классификатор может получить только в том случае, если дает верные прогнозы чаще, чем по такой упрощенной стратегии.

Значения каппа-статистики лежат в диапазоне от 0 до 1. Единица указывает на идеальное соответствие прогнозов модели и истинных значений. Значения меньше единицы указывают на неидеальное соответствие. В зависимости от того, как будет использоваться модель, возможны различные варианты интерпретации каппа-статистики. Вот один из типичных вариантов:

- плохое соответствие — менее 0,20;
- удовлетворительное соответствие — от 0,20 до 0,40;
- среднее соответствие — от 0,40 до 0,60;
- хорошее соответствие — 0,60 до 0,80;

очень хорошее соответствие — от 0,80 до 1,00.

Важно отметить, что эти категории являются субъективными. С одной стороны, «хорошее соответствие» может быть более чем достаточным для прогнозирования любимого вкуса мороженого, а с другой — «очень хорошего соответствия» может оказаться недостаточно, если ваша цель — выявить врожденные дефекты.



Подробнее об этом читайте в статье: Landis J.R., Koch G.G. The measurement of observer agreement for categorical data. *Biometrics*, 1997. Vol. 33. P. 159–174.

Ниже приведена формула для вычисления каппа-статистики. В этой формуле $\Pr(a)$ означает пропорцию реального соответствия, а $\Pr(e)$ — ожидаемое соответствие между прогнозом классификатора и истинными значениями при условии, что они выбраны случайным образом:

$$\kappa = \frac{\Pr(a) - \Pr(e)}{1 - \Pr(e)} .$$



Есть несколько способов вычисления каппа-статистики. Самый распространенный метод, представленный здесь, использует коэффициент каппа Коэна, описанный в статье: Cohen J. A coefficient of agreement for nominal scales // *Education and Psychological Measurement*, 1960. Vol. 20. P. 37–46.

Эти пропорции легко получить из матрицы несоответствий, если знать, где искать. Рассмотрим матрицу несоответствий для модели классификации СМС, построенную с помощью функции `CrossTable()`, которая для удобства приводится еще раз:

sms_results\$actual_type	sms_results\$predict_type		Row Total
	ham	spam	
ham	1203 16.128 0.997 0.975 0.865	4 127.580 0.003 0.026 0.003	1207 0.868
spam	31 106.377 0.169 0.025 0.022	152 841.470 0.831 0.974 0.109	183 0.132
Column Total	1234 0.888	156 0.112	1390

Нижнее значение в каждой ячейке означает долю всех экземпляров, попадающих в эту ячейку. Поэтому, чтобы вычислить наблюдаемое

соответствие $Pr(a)$, достаточно добавить пропорцию всех случаев, когда прогнозируемый и реальный типы СМС совпадают. Таким образом, можно вычислить $Pr(a)$ следующим образом:

```
> pr_a <- 0.865 + 0.109> pr_a[1] 0.974
```

Для этого классификатора наблюдаемые и фактические значения совпадают в 97,4 % случаев — легко заметить, что это соответствует точности. Каппа-статистика корректирует точность относительно ожидаемого соответствия $Pr(e)$, которое представляет собой вероятность того, что совпадение спрогнозированных и реальных значений может оказаться случайным, исходя из предположения, что оба они были выбраны случайным образом, в соответствии с наблюдаемыми пропорциями.

Чтобы получить эти наблюдаемые пропорции, можно использовать правила вычисления вероятности, которые рассмотрены в главе 4. Если предположить, что эти два события независимы (то есть одно не влияет на другое), то согласно правилам вычисления вероятности вероятность возникновения обоих событий равна произведению вероятностей возникновения каждого из них. Например, мы знаем, что вероятность того, что письмо не является спамом, равна:

$Pr(\text{actual_type не спам}) \times Pr(\text{predicted_type не спам})$.

А вероятность того, что оба являются спамом:

$Pr(\text{actual_type — спам}) \times Pr(\text{predicted_type — спам})$.

Вероятность того, что прогнозируемый или реальный тип является спамом либо обычным письмом, может быть вычислена на основании суммы строки или столбца. Например, $Pr(\text{actual_type не спам}) = 0,868$, а $Pr(\text{predicted type не спам}) = 0,888$.

$Pr(e)$ можно вычислить как сумму вероятностей того, что прогнозируемые и реальные значения совпадают, будь сообщение спамом или нет. Напомню, что для взаимоисключающих событий (то есть таких, которые не могут происходить одновременно) вероятность одного из событий равна сумме их вероятностей. Поэтому, чтобы получить итоговое значение $Pr(e)$, достаточно сложить оба произведения:

```
> pr_e <- 0.868 * 0.888 + 0.132 * 0.112>
pr_e[1] 0.785568
```

Поскольку $Pr(e)$ равно 0,786, то вероятность случайного совпадения прогнозируемых и реальных значений составляет примерно 78,6 % случаев.

Это означает, что теперь у нас есть вся информация, необходимая для окончательного вычисления формулы каппа-статистики. Подставляя в эту формулу значения $Pr(a)$ и $Pr(e)$, получим:

```
> k <- (pr_a - pr_e) / (1 - pr_e)> k[1]
0.8787494
```

Коэффициент равен примерно 0,88, что соответствует представленному выше результату `confusionMatrix()` из пакета `caret` (небольшое расхождение вызвано округлением). Используя предложенную интерпретацию, отмечу, что существует очень хорошее соответствие между прогнозами классификатора и реальными значениями.

В R есть пара функций для автоматического расчета коэффициентов каппа. В функции `Kappa()` (обратите внимание: с заглавной буквой K) из пакета `Visualizing Categorical Data (vcd)` используется матрица несоответствий прогнозируемых и реальных значений. Чтобы вычислить значение каппа, нужно установить пакет, ввести команду `install.packages("vcd")` и затем использовать следующие команды:

```
> Kappa(table(sms_results$actual_type,
sms_results$predict_type))
value
ASEUnweighted 0.8825203 0.01949315Weighted
0.8825203 0.01949315
```

Нас интересует невзвешенная каппа-статистика. Значение 0,88 соответствует тому, что мы вычислили вручную.

Взвешенная каппа-статистика используется в тех случаях, когда есть различные степени соответствия. Например, при наличии шкалы с градациями «холодно», «прохладно», «тепло» и «горячо» значение «тепло» больше соответствует варианту «горячо», чем варианту «холодно». Для события с двумя исходами, такими как «спам» и «не спам», взвешенная и невзвешенная каппа-статистика будут идентичны.

Функция `kappa2()` из пакета `Interrater Reliability (irr)` позволяет вычислять каппа-статистику на основе векторов прогнозируемых и реальных значений, хранящихся во фрейме данных. Для получения каппа-статистики с помощью этой функции нужно установить пакет с помощью команды `install.packages("irr")` и ввести следующие команды:

```
> kappa2(sms_results[1:2])Cohen's Kappa for 2
Raters (Weights: unweighted)Subjects =
1390 Raters = 2 Kappa = 0.883 z = 33p-
value = 0
```

Функции `Kappa()` и `kappa2()` вычисляют одно и то же значение каппа-статистики, поэтому вы можете использовать любой вариант, который сочтете удобным.



Однако будьте внимательны: не используйте встроенную функцию `kappa()`. Она не имеет никакого отношения к описанной здесь каппа-статистике!

Чувствительность и специфичность

Подобрать полезный классификатор часто означает найти баланс между чрезмерно консервативными и чрезмерно агрессивными прогнозами. Например, фильтр электронной почты может гарантированно устранять все спам-сообщения, одновременно агрессивно удаляя почти всю обычную почту. И наоборот, для того, чтобы гарантировать, что ни одно полезное сообщение не будет случайно отфильтровано, фильтр может пропускать неприемлемо большое количество спама. Этот компромисс отражает следующая пара показателей эффективности: чувствительность и специфичность.

Чувствительность (sensitivity) модели (также называемая *частотой истинно положительных прогнозов*) — это доля правильно классифицированных положительных примеров. Поэтому, как показано в следующей формуле, чувствительность вычисляется как количество истинно положительных прогнозов, разделенное на общее количество положительных результатов, классифицированных как правильно (истинно положительные), так и ошибочно (ложноотрицательные):

$$\text{sensitivity} = \frac{TP}{TP + FN}.$$

Специфичность (specificity) модели (или *частота истинно отрицательных прогнозов*) — это доля правильно классифицированных отрицательных примеров. Подобно чувствительности, специфичность вычисляется как число истинно отрицательных прогнозов, разделенное на общее количество отрицательных результатов — как истинно отрицательных, так и ложноположительных.

$$\text{specificity} = \frac{TN}{TN + FP}.$$

На основании матрицы несоответствий для классификатора СМС можно легко вычислить эти показатели вручную. Предполагая, что спам является положительным классом, мы можем убедиться, что числа в результатах `confusionMatrix()` верны. Например, вот расчет чувствительности:

```
> sens <- 152 / (152 + 31) > sens[1] 0.8306011
```

Аналогично специфичность вычисляется следующим образом:

```
> spec <- 1203 / (1203 + 4) > spec[1] 0.996686
```

В пакет `caret` входят функции для вычисления чувствительности и специфичности непосредственно на основе векторов прогнозируемых и реальных значений. Используя эти функции, убедитесь, что правильно указываете параметры `positive` и `negative`:

```
> library(caret)
sensitivity(sms_results$predict_type,
sms_results$actual_type,           positive =
```

```
"spam" ) [1] 0.8306011 >  
specificity(sms_results$predict_type,  
sms_results$actual_type, negative =  
"ham" ) [1] 0.996686
```

Чувствительность и специфичность принимают значения от 0 до 1, при этом значения, близкие к единице, являются более желательными. Конечно, важно найти баланс между этими двумя показателями — как правило, эта задача очень сильно зависит от предметной области.

В данном примере чувствительность 0,831 означает, что 83,1 % спамовых сообщений были классифицированы правильно. Аналогичным образом специфичность, равная 0,997, подразумевает, что 99,7 % сообщений, не являющихся спамом, были правильно классифицированы. Другими словами, 0,3 % полезных сообщений были отклонены как спам. Идея отклонения 0,3 % полезных СМС может оказаться неприемлемой, а может быть сочтена разумным компромиссом — при условии уменьшения количества спама.

Чувствительность и специфичность — это инструменты для поиска таких компромиссов. Обычно изменение модели и тестирование различных моделей продолжается до тех пор, пока не будет найдена такая модель, которая бы соответствовала требуемому порогу чувствительности и специфичности. Визуализации, подобные тем, что будут показаны далее в этой главе, также помогают понять, где находится баланс между чувствительностью и специфичностью.

Точность и полнота

С чувствительностью и специфичностью тесно связаны две другие характеристики, определяющие компромиссы в классификации: точность и полнота. Эти статистические показатели используются в основном в задачах поиска информации и предназначены для того, чтобы показать, насколько интересны и актуальны результаты модели или же прогнозы разбавлены бессмысленным шумом.

Точность (precision) (или *прогностическая ценность положительных результатов*) определяется как доля реальных положительных примеров, которые спрогнозированы как положительные. Другими словами, как часто модель правильно прогнозирует положительный класс? Точная модель будет прогнозировать положительный класс только в тех случаях, которые действительно могут быть положительными. Это будет весьма заслуживающая доверия модель.

Как вы думаете, что произойдет, если модель будет очень неточной? Со временем ее результатам будут все меньше доверять. В контексте поиска информации это будет похоже на поисковую систему, такую как Google, которая возвращает нерелевантные результаты. В конце концов

пользователи уйдут к конкуренту, например к Bing. В случае спам-фильтра СМС высокая точность означает, что модель способна точно фильтровать только спам, пропуская полезные сообщения.

$$\text{precision} = \frac{TP}{TP + FP}.$$

Полнота (recall), напротив, показатель того, насколько полными являются результаты. Как показано в следующей формуле, полнота определяется как доля истинно положительных прогнозов в общем количестве положительных прогнозов. Возможно, вы посчитали это чувствительностью; однако интерпретация немного отличается.

Модель с высоким уровнем полноты фиксирует большую часть положительных примеров, а это означает, что она имеет широкий охват. Например, поисковая система с высокой полнотой возвращает большое количество документов, отвечающих поисковому запросу. Аналогично спам-фильтр СМС имеет высокий уровень полноты, если правильно идентифицирует большинство спам-сообщений.

$$\text{recall} = \frac{TP}{TP + FN}.$$

Точность и полноту можно рассчитать на основании матрицы несоответствий. Если, как и раньше, предположить, что спам является положительным классом, то точность вычисляется так:

```
> prec <- 152 / (152 + 4) > prec[1] 0.974359
```

А полнота — так:

```
> rec <- 152 / (152 + 31) > rec[1] 0.8306011
```

Пакет `caret` позволяет вычислить любой из этих показателей на основе векторов спрогнозированных и реальных классов. Для вычисления точности используется функция `posPredValue()`:

```
> library(caret) >
posPredValue(sms_results$predict_type,
sms_results$actual_type, positive =
"spam") [1] 0.974359
```

Для вычисления полноты применяется уже знакомая нам функция `sensitivity()`:

```
> sensitivity(sms_results$predict_type,
sms_results$actual_type, positive =
"spam") [1] 0.8306011
```

Подобно внутреннему компромиссу между чувствительностью и специфичностью, для большинства реальных задач трудно построить модель, которая бы обладала одновременно и высокой точностью, и высокой полнотой. Легко быть точным, если собирать с дерева только низко висящие фрукты — примеры, которые легко классифицировать. Аналогичным образом легко создать модель с высокой полнотой, построив

очень широкую сеть, — это будет означать, что модель слишком агрессивна при выявлении положительных случаев. И наоборот, очень трудно обеспечить одновременно и высокую точность, и высокую полноту. Поэтому важно протестировать множество моделей, чтобы найти такое сочетание точности и полноты, которое бы отвечало потребностям вашего проекта.

F-мера

Мера эффективности модели, которая объединяет точность и полноту в одно число, называется *F-мерой* (F-measure) (также называемой *оценкой F1* или *F-оценкой*). F-мера сочетает в себе точность и полноту, используя среднее гармоническое — тип среднего, который описывает скорость изменения величины. Используется именно среднее гармоническое, а не более привычное среднее арифметическое, поскольку точность и полнота описываются как пропорции, принимающие значения от 0 до 1, и поэтому их можно интерпретировать как скорости. F-мера вычисляется по следующей формуле:

$$F\text{-measure} = \frac{2 \times \text{precision} \times \text{recall}}{\text{recall} + \text{precision}} = \frac{2 \times TP}{2 \times TP + FP + FN}.$$

Чтобы вычислить F-меру, можно воспользоваться рассчитанными ранее значениями точности и полноты:

```
> f <- (2 * prec * rec) / (prec + rec) > f[1]
0.8967552
```

Получится точно то же самое, как при использовании расчетов по матрице несоответствий:

```
> f <- (2 * 152) / (2 * 152 + 4 + 31) > f[1]
0.8967552
```

Поскольку F-мера описывает эффективность модели одним числом, она является удобным способом сравнения нескольких моделей между собой. Однако это предполагает, что точность и полнота имеют равный вес, что не всегда верно. Можно вычислить F-меру, используя разный вес для точности и полноты, но выбор весов может быть сложным в лучшем случае и произвольным в худшем. Рекомендуется использовать такие показатели, как F-мера, в сочетании с методами, которые учитывают достоинства и недостатки модели в более глобальном масштабе: например, так, как описано в следующем разделе.

Визуализация компромиссов эффективности с помощью ROC-кривых

Визуализация помогает лучше исследовать эффективность алгоритмов машинного обучения. Если статистические данные, такие как чувствительность и специфичность или точность и полнота, пытаются

свести эффективность модели к одному числу, то визуализации показывают, как алгоритм обучения работает в широком диапазоне условий.

Поскольку у каждого алгоритма обучения свои особенности, может оказаться так, что две модели с одинаковой точностью будут резко отличаться тем, как именно они достигают такой точности. Некоторые модели могут с трудом делать определенные типы прогнозов, в то время как другие с легкостью справляются с такими прогнозами, но при этом пропускают случаи, с которыми третьи алгоритмы справляются легко. Визуализация является способом увидеть эти компромиссы, сравнивая алгоритмы обучения между собой на одном графике.

Для изучения компромисса между обнаружением истинно положительных результатов, избегая при этом ложноположительных результатов, обычно используется *кривая рабочих характеристик приемника* (Receiver Operating Characteristic, ROC). Как несложно догадаться из названия, ROC-кривые изначально разрабатывались инженерами-связистами. Во время Второй мировой войны операторы радаров и радисты использовали ROC-кривые для измерения способности приемника различать истинные и ложные сигналы тревоги. Сегодня подобная методика применяется для визуализации эффективности моделей машинного обучения.

На рис. 10.4 показаны характеристики типичной ROC-диаграммы. По вертикальной оси откладывается частота истинно положительных прогнозов, а по горизонтальной — частота ложноположительных прогнозов. Поскольку эти значения эквивалентны чувствительности и (1-специфичности) соответственно, диаграмма также называется графиком чувствительности/специфичности.

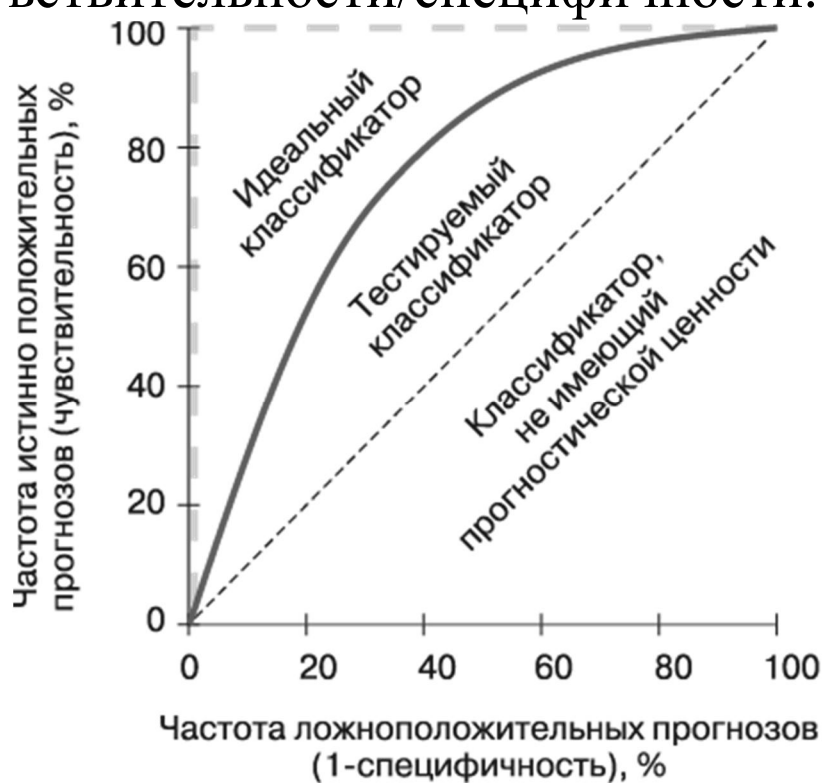


Рис. 10.4. ROC-кривая отображает форму тестируемого классификатора относительно идеального и бесполезного классификаторов

Точки, образующие ROC-кривую, показывают частоту истинно положительных прогнозов при различных порогах ложноположительных прогнозов. Для построения таких кривых прогнозы классификатора сортируются по оцениваемой модели вероятности положительной классификации, причем сначала идут самые большие значения. Начиная с исходной точки, влияние каждого прогноза на частоту истинно положительных и ложноположительных прогнозов приводит к тому, что кривая будет отклоняться ближе к вертикали (в случае правильного прогноза) или к горизонтали (при неверном прогнозе).

Чтобы объяснить эту концепцию, на рис. 10.4 представлены три гипотетических классификатора: диагональная линия, идущая от нижнего левого к верхнему правому углу диаграммы, представляет *классификатор, не имеющий прогностической ценности*. Такой классификатор с одинаковой частотой дает истинно положительные и ложноположительные прогнозы, что означает, что такой классификатор не способен отличить одно от другого. Это базовый уровень, по которому можно судить о других классификаторах. ROC-кривые, расположенные близко к этой линии, указывают на то, что соответствующие модели не особенно полезны. И наоборот, кривая *идеального классификатора* проходит через точку со 100%-ной частотой истинно положительных прогнозов и нулевой частотой ложноположительных прогнозов. Такой классификатор правильно идентифицирует все истинно положительные результаты, прежде чем неправильно классифицировать любой отрицательный результат. Большинство реальных классификаторов схожи с тестовым в том, что находятся где-то посередине между идеальным и бесполезным классификаторами.

Чем ближе кривая к идеальному классификатору, тем он лучше определяет положительные значения. Это можно измерить с помощью статистического показателя, известного как *площадь под ROC-кривой* (Area Under the ROC Curve, AUC). AUC рассматривает ROC-диаграмму как двумерный квадрат и измеряет общую площадь под ROC-кривой. AUC лежит в пределах от 0,5 (для классификатора, не имеющего прогностической ценности) до 1,0 (для идеального классификатора). Существует следующее соглашение для интерпретации баллов AUC, аналогично системе оценки успеваемости в США:

- **A** (превосходно) — от 0,9 до 1,0;
- **B** (отлично/хорошо) — от 0,8 до 0,9;
- **C** (приемлемо/удовлетворительно) — от 0,7 до 0,8;
- **D** (плохо) — от 0,6 до 0,7;
- **E** (без оценки) — от 0,5 до 0,6.

Как и в большинстве подобных шкал, для одних задач эти уровни могут подходить лучше, чем для других; эта классификация несколько субъективна.

Как видно на рис. 10.5, две ROC-кривые могут иметь очень разную форму и при этом AUC будет идентична. Поэтому одной лишь AUC недостаточно для определения «лучшей» модели. Самый надежный метод — использовать AUC в сочетании с качественным исследованием ROC-кривой.

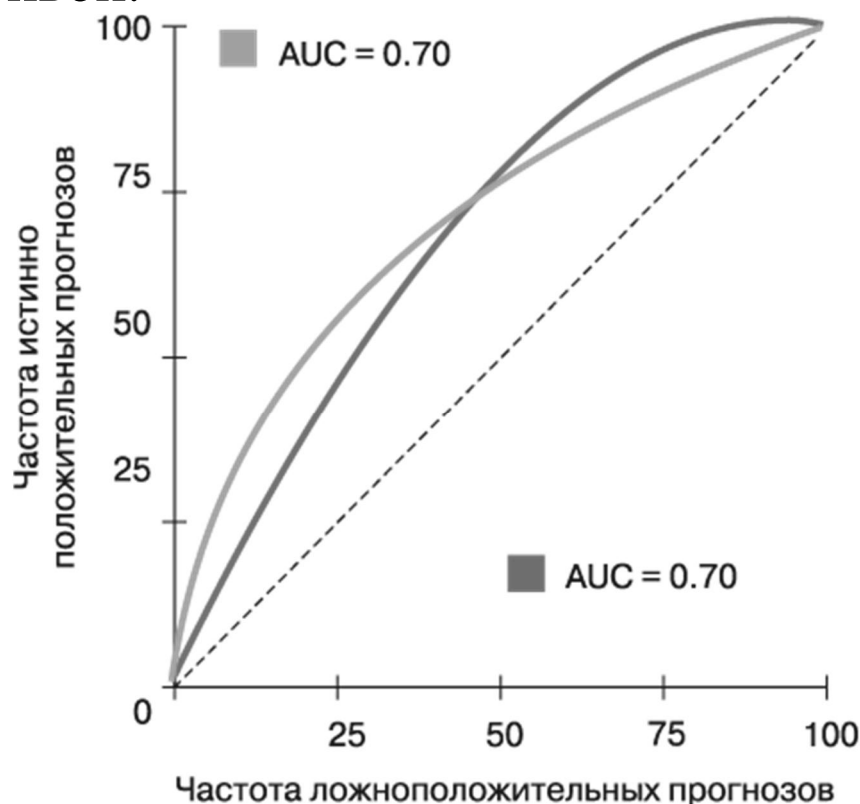


Рис. 10.5. Несмотря на одинаковую AUC ROC-кривых, модели могут иметь разную эффективность

Пакет `pROC` предоставляет удобный набор функций для построения ROC-кривых и вычисления AUC. На сайте <https://web.expasy.org/pROC/> вы найдете полный список функций, а также несколько примеров визуализации. Прежде чем продолжить работу, не забудьте установить этот пакет с помощью команды `install.packages("pROC")`.



Подробнее о пакете `pROC` читайте в статье: Robin X., Turck N., Hainard A., Tiberti N., Lisacek F., Sanchez J.C., Mueller M. `pROC`: an open-source package for R and S+ to analyze and compare ROC curves *BMC. Bioinformatics*, 2011. P. 12–77.

Для построения визуализаций с помощью пакета `pROC` необходимы два вектора данных: один содержит предполагаемую вероятность прогнозирования положительного класса, а второй — прогнозируемые значения класса. В случае классификатора СМС оценочные вероятности классификации писем как спам и реальные метки классов для функции `roc()` предоставляются следующим образом:

```
> library(pROC) > sms_roc <-  
roc(sms_results$prob_spam,  
sms_results$actual_type)
```

Используя объект `sms_roc`, можно построить ROC-кривую с помощью функции `plot()`. Как показано в следующем фрагменте кода, у этой функции есть множество стандартных параметров для настройки визуализации, таких как `main` (заголовок), `col` (цвет линии) и `lwd` (ширина линии). Параметр `legacy.axes` в пакете `pROC` позволяет откладывать по оси X значения (1-специфичность), что является распространенным соглашением:

```
> plot(sms_roc, main = "ROC curve for SMS spam  
filter", col = "blue", lwd = 2, legacy.axes  
= TRUE)
```

Конечный результат представлен на рис. 10.6 — это ROC-кривая с диагональной контрольной линией, соответствующей базовому классификатору без прогностической ценности.

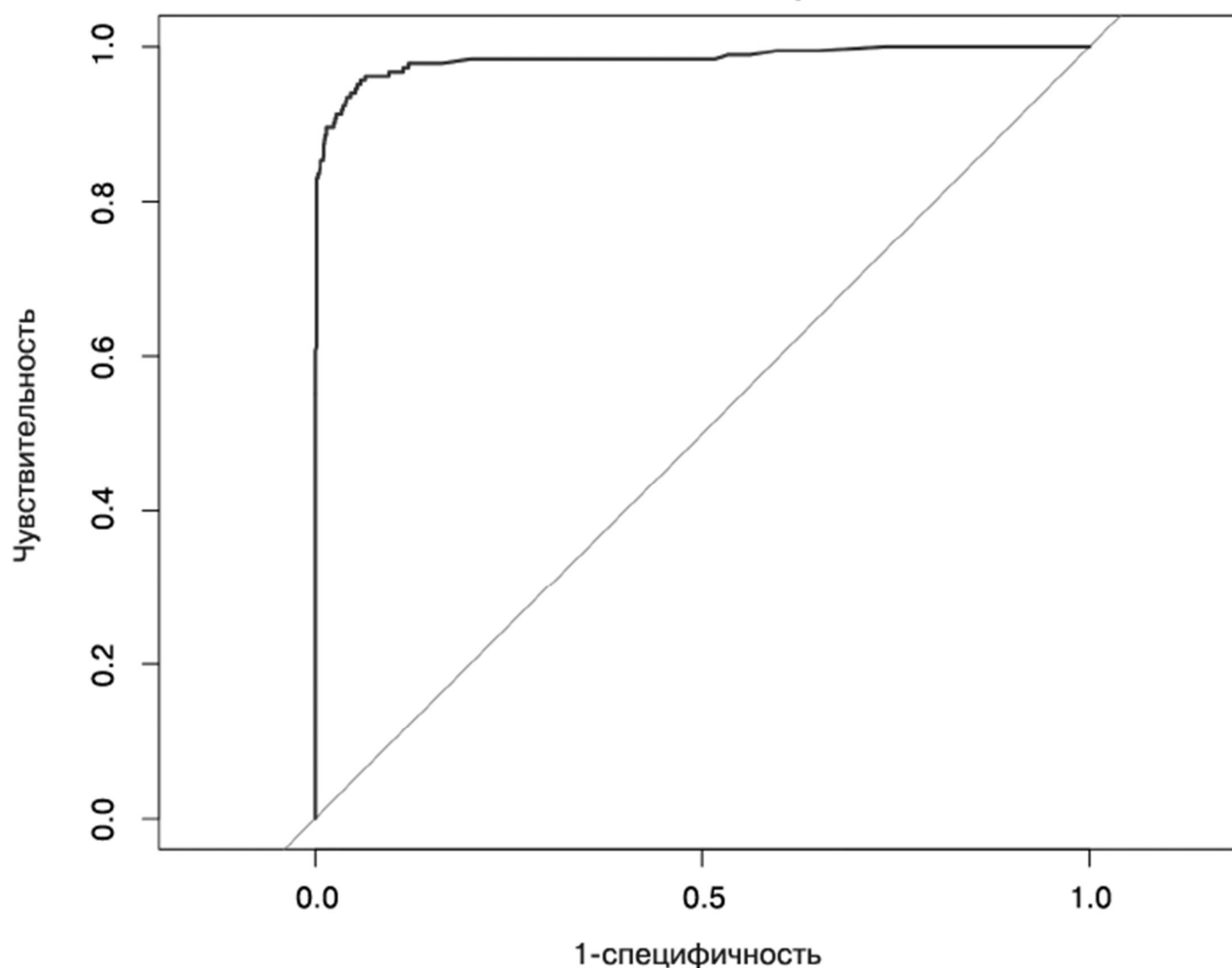


Рис. 10.6. ROC-кривая для СМС-классификатора, использующего наивный байесовский алгоритм

Сразу видно, что ROC-кривая занимает пространство в верхнем левом углу диаграммы. Это говорит о том, что она ближе к идеальному классификатору, чем пунктирная линия, соответствующая бесполезному классификатору.

Чтобы сравнить эффективность этой модели с другими моделями, делающими прогнозы на том же наборе данных, можно добавить на этот

же график дополнительные ROC-кривые. Предположим, что мы также обучили на данных СМС модель k-NN с применением функции `knn()`, описанной в главе 3. С помощью этой модели были вычислены прогнозируемые вероятности спама для каждой записи в наборе тестов. Затем они были сохранены в CSV-файле, который доступен для загрузки. Загрузив файл, для построения ROC-кривой, как и прежде, применим функцию `roc()`, а затем воспользуемся функцией `plot()` с параметром `add=TRUE`, чтобы разместить эту кривую на предыдущем графике:

```
> sms_results_knn <-  
read.csv("sms_results_knn.csv")> sms_roc_knn <-  
roc(sms_results$actual_type,  
sms_results_knn$p_spam)> plot(sms_roc_knn, col =  
"red", lwd = 2, add = TRUE)
```

На рис. 10.7 видно, что появилась вторая кривая, соответствующая эффективности модели k-NN, которая делает прогнозы на том же наборе тестов, что и модель, построенная по наивному байесовскому алгоритму. Кривая для модели k-NN в каждой точке размещается ниже, чем первая модель. Это позволяет предположить, что данная модель неизменно хуже, чем наивный байесовский алгоритм.

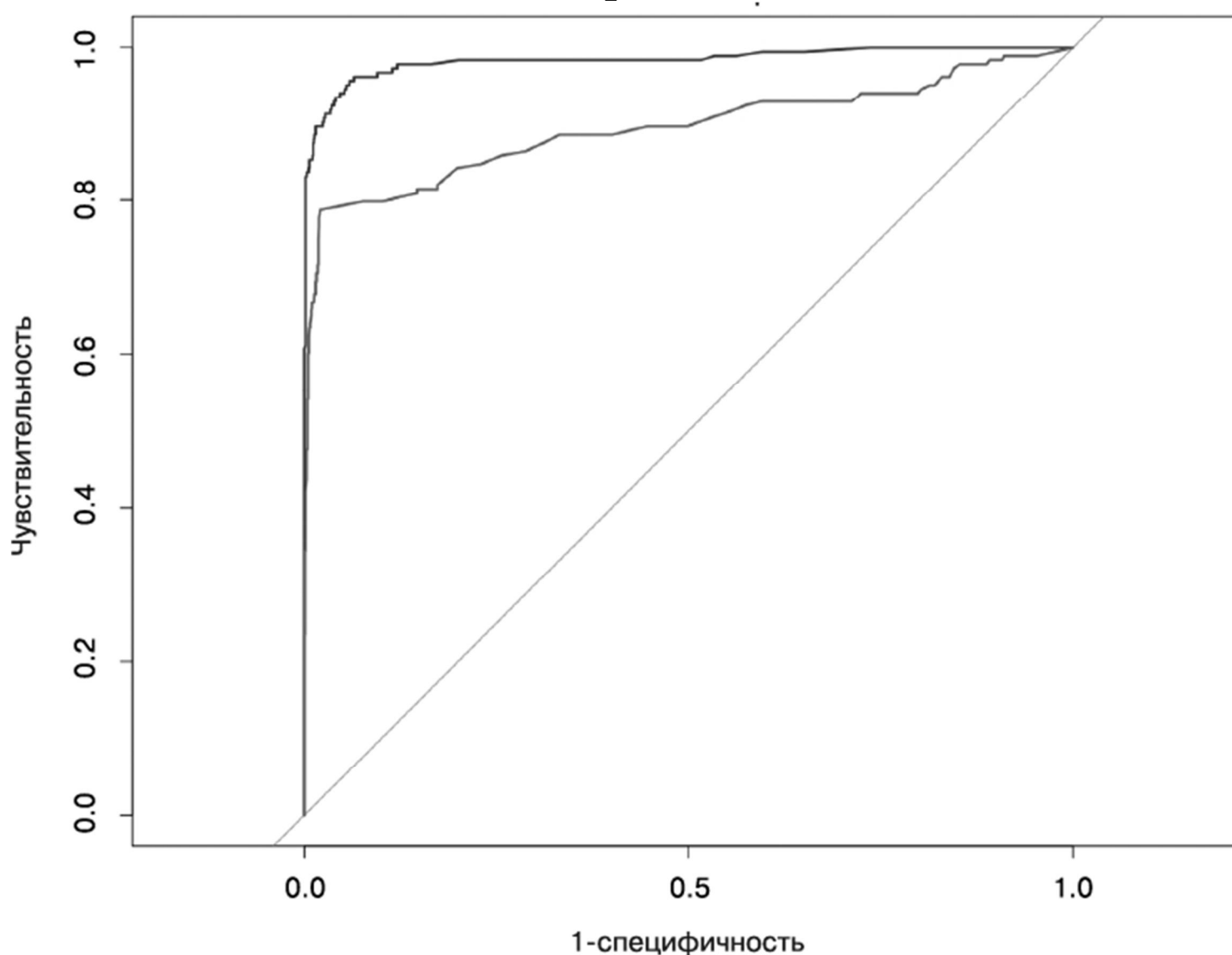


Рис. 10.7. ROC-кривые, позволяющие сравнить эффективность наивного байесовского алгоритма (верхняя кривая) и алгоритма k-NN на тестовом наборе СМС

Чтобы подтвердить это количественно, можно воспользоваться функциями пакета `rROC` для расчета AUC. Для этого мы применим

функцию `auc()` к объекту `sms_roc` каждой модели, как показано в следующем коде:

```
> auc(sms_roc)Area under the curve: 0.9836>
auc(sms_roc_knn)Area under the curve: 0.8942
```

AUC для СМС-классификатора, построенного по наивному байесовскому алгоритму, составляет 0,98, что чрезвычайно высоко и существенно лучше, чем AUC классификатора k-NN, которая равна 0,89. Но как узнать, будет ли эта модель столь же хорошо работать на другом наборе данных, или же различие оказалось больше ожидаемого лишь по случайности? Для того чтобы ответить на вопрос, нужно лучше понять, насколько далеко можно экстраполировать прогнозы модели за пределы тестовых данных.



Как уже отмечалось, одного лишь значения AUC часто бывает недостаточно для определения лучшей модели. В данном примере AUC идентифицирует лучшую модель, потому что ROC-кривые не пересекаются. В других случаях лучшая модель будет зависеть от того, как именно эта модель будет использоваться. Если ROC-кривые пересекаются, их можно объединить в еще более сильную модель, используя методы, описанные в главе 11.

Оценка эффективности в будущем

В некоторых R-пакетах ML представлены матрицы несоответствий и показатели эффективности, вычисляемые в процессе построения модели. Цель этих статистических показателей — дать представление об *эмпирическом риске* модели, который возникает, когда данные обучения спрогнозированы неверно, несмотря на то что модель построена непосредственно на этих данных. Такая информация может использоваться в качестве грубой диагностики для выявления явно плохих показателей.

Однако для оценки будущей эффективности эмпирический риск не особенно полезен. Например, модель, в которой для точной классификации каждого обучающего экземпляра используется механическое запоминание, имеет нулевой эмпирический риск, однако эта модель не может обобщить свои прогнозы на данные, которые ей никогда ранее не встречались. Поэтому частота ошибок в тренировочных данных может быть чересчур оптимистичной оценкой в отношении будущих характеристик модели.

Вместо того чтобы полагаться на эмпирический риск, лучше оценить эффективность модели на данных, которые ей еще не встречались. Такой подход применялся в предыдущих главах, когда доступные данные были разделены на два набора: для обучения и для тестирования. Однако бывают случаи, когда создавать тренировочные и тестовые наборы данных

не очень удобно. Например, в ситуации, когда есть только небольшой пул данных, не все захотят уменьшать выборку еще больше.

К счастью, существуют другие способы оценки эффективности модели на новых данных. В состав пакета `caret`, который мы уже использовали для вычисления показателей эффективности, также входит ряд функций для оценки будущих результатов. Если вы лишь просматривали примеры R-кода, но еще не установили пакет `caret`, сделайте это. Вам также потребуется загрузить пакет в R-сессию с помощью команды `library(caret)`.

Метод отложенных данных

Процедура разделения данных на тренировочный и тестовый наборы, к которой мы прибегали в предыдущих главах, называется *методом отложенных данных*. На рис. 10.8 показано, как *тренировочный набор данных* используется для формирования модели, которая затем применяется к *тестовому набору данных* для генерации прогнозов, которые далее подвергаются оценке. Как правило, примерно треть данных оставляется для тестирования, а две трети используются для обучения, но эта пропорция может быть разной, в зависимости от объема доступных данных. Чтобы между тренировочным и тестовым наборами данных не оказалось систематических различий, примеры данных разделяются на эти две группы случайным образом.

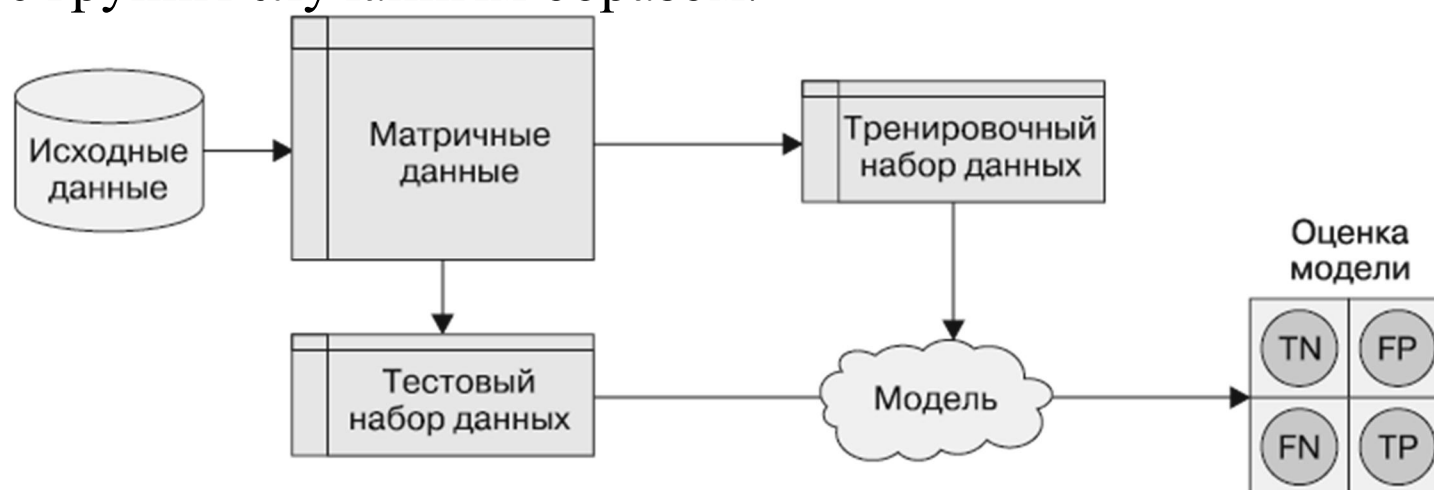


Рис. 10.8. Простейший метод отложенных данных заключается в разделении данных на тренировочный и тестовый наборы

Для того чтобы метод отложенных данных позволил получить действительно точную оценку будущей эффективности, ни в коем случае нельзя допускать, чтобы эффективность на тестовом наборе данных оказывала влияние на модель. Это правило легко нарушить неосознанно, выбрав лучшую модель по результатам повторного тестирования. Предположим, что мы построили на тренировочных данных несколько моделей и выбрали из них ту, которая обеспечивает самую высокую точность на тестовых данных. Поскольку мы выбрали модель с наилучшим результатом, эффективность на тестовом наборе не гарантирует, что модель будет такой же эффективной на новых данных.



Внимательный читатель заметит, что в предыдущих главах тестовые данные использовались как для оценки модели, так и для повышения ее эффективности. Это было сделано в иллюстративных целях: на самом деле это нарушение изложенного выше правила. Следовательно, продемонстрированная статистика эффективности модели не была достоверной оценкой ее эффективности в будущем, на новых данных. Чтобы избежать такой проблемы, лучше разделить исходные данные так, чтобы, помимо тренировочного и тестового наборов данных, был еще *валидационный набор данных*. Его можно использовать для итерационного улучшения выбранной модели или нескольких моделей, так что тестовый набор данных останется только для однократного применения в виде завершающего этапа для определения предполагаемой частоты ошибок в будущих прогнозах (рис. 10.9). Как правило, исходные данные разделяются на тренировочный, тестовый и валидационный наборы в соотношении 50, 25 и 25 % соответственно.

В простейшем случае при создании наборов отложенных данных для назначения записей в тот или иной раздел используются генераторы случайных чисел. Такой метод был впервые использован в главе 5 для создания обучающего и тестового наборов данных.



Для того чтобы выполнить приведенные ниже примеры, скачайте набор данных `credit.csv` и загрузите его во фрейм данных с помощью команды `credit <- read.csv("credit.csv")`.

Предположим, у нас есть фрейм данных с названием `credit` и 1000 строк данных. Можно разбить его на три раздела следующим образом. Сначала создадим вектор идентификаторов строк от 1 до 1000, расположенных в произвольном порядке. Для этого воспользуемся функцией `runif()`, которая по умолчанию генерирует заданное количество случайных значений в диапазоне от 0 до 1. Функция `runif()` получила такое название от случайного равномерного (random uniform) распределения, о котором говорилось в главе 2.

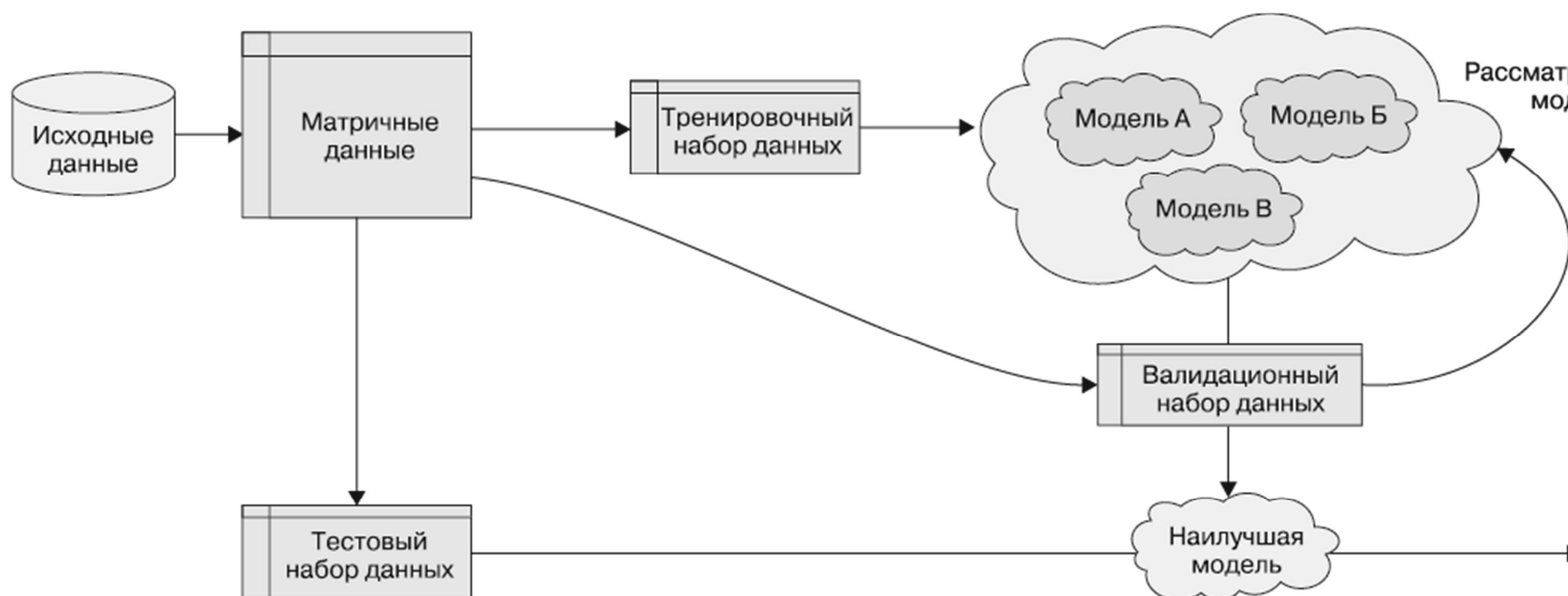


Рис. 10.9. Валидационный набор данных может быть исключен из обучения для выбора из нескольких моделей

Затем функция `order()` возвращает вектор, содержащий последовательность упорядоченных 1000 случайных чисел. Например, `order(c(0.5, 0.25, 0.75, 0.1))` возвращает последовательность `4213`: наименьшее из чисел (0,1) отображается четвертым, предпоследнее по величине (0,25) — вторым и т.д.

```
> random_ids <- order(runif(1000))
```

Затем фрейм данных `credit` делится на группы из 500, 250 и 250 записей, соответствующих тренировочному, тестовому и валидационному наборам данных. Идентификаторы примеров, попадающих в каждый из этих наборов, выбираются случайным образом:

```
> credit_train <- credit[random_ids[1:500], ]>
credit_validate <- credit[random_ids[501:750], ]>
credit_test <- credit[random_ids[751:1000], ]
```

Одна из проблем отложенных данных состоит в том, что в каждый набор данных может попасть одних классов больше, чем других. Если один или несколько классов составляют очень малую часть набора данных, это может привести к тому, что классы будут исключены из тренировочного набора данных, что является серьезной проблемой, поскольку тогда модель не сможет изучить эти классы.

Чтобы уменьшить вероятность этого, можно использовать метод *стратифицированной случайной выборки*. Случайная выборка, как правило, содержит значения всех классов примерно в тех же пропорциях, что и полный набор данных, однако стратифицированная случайная выборка гарантирует, что в каждом наборе все классы будут представлены почти в тех же пропорциях, что и в полном наборе, даже если некоторые из них небольшие.

В пакет `caret` входит функция `createDataPartition()`, которая создает разделы набора данных на основе стратифицированных отложенных данных. Далее приведены команды для создания стратифицированной выборки тренировочных и тестовых данных для набора данных `credit`. Аргументами функции `createDataPartition()` являются вектор значений класса (в этом случае `default` содержит данные о том, была ли неуплата долгов по кредиту), а также параметр `p`, который указывает на то, какая доля экземпляров должна быть включена в раздел. Параметр `list=FALSE` не дает сохранять результаты в виде объекта списка:

```
> in_train <-  
createDataPartition(credit$default, p = 0.75,  
list = FALSE)> credit_train <- credit[in_train,  
> credit_test <- credit[-in_train, ]
```

Вектор `in_train` содержит номера строк, включенных в обучающую выборку. Эти номера строк можно использовать, чтобы выбрать примеры для фрейма данных `credit_train`. Аналогично, добавив знак «минус», можно внести строки, отсутствующие в векторе `in_train`, для набора данных `credit_test`.

При этом классы распределяются равномерно, однако стратифицированная выборка не гарантирует другие типы репрезентативности. В выборках может оказаться слишком много или слишком мало сложных случаев, а также легко прогнозируемых случаев или выбросов. Это особенно верно для небольших наборов данных, в которых может не быть достаточного количества таких случаев, чтобы их можно было разделить между тренировочными и тестовыми наборами.

Кроме потенциально искаженных выборок, с методом отложенных данных связана еще одна проблема. Она состоит в том, что значительная часть данных должна быть зарезервирована для тестирования и валидации модели. Поскольку эти данные нельзя использовать для обучения модели до тех пор, пока не будут измерены ее характеристики, оценки эффективности, вероятно, будут слишком консервативными.



Поскольку модели, обученные на больших наборах данных, как правило, работают лучше, чем модели, обученные на меньшем количестве данных, часто выполняется повторное обучение модели на полном наборе данных (то есть на наборе, который включает в себя тренировочный и тестовый наборы, а также валидационный набор) после того, как была выбрана окончательная модель и выполнена оценка ее эффективности.

Для разрешения проблемы произвольно составленных тренировочных наборов данных иногда используется метод *повторно отложенных данных*. Это частный случай метода отложенных данных, в котором для оценки эффективности модели применяется усредненный результат, составленный из нескольких случайных выборок отложенных данных. Поскольку используются несколько выборок отложенных данных, это снижает вероятность обучения или тестирования модели на нерепрезентативных данных. Об этом — в следующем разделе.

Кросс-валидация

Метод повторно отложенных данных является основой метода *кросс-валидации* (Cross-Validation, CV) *по k-блокам*. Этот метод стал отраслевым стандартом для оценки эффективности модели. Вместо того чтобы делать повторные случайные выборки, которые потенциально могут использовать одну и ту же запись несколько раз, CV по k-блокам случайным образом делит данные на k совершенно разных случайных разделов — *блоков*.

Параметру *k* можно присвоить любое числовое значение, однако чаще всего используется CV по десяти блокам. Почему блоков именно десять? Как показывают эмпирические данные, использование большего числа блоков дает мало дополнительных преимуществ. Для каждого из десяти блоков (каждый из которых содержит 10 % общего объема данных), модель машинного обучения строится на оставшихся 90 % данных. Затем 10 % примеров, содержащихся в блоке, используется для оценки модели. После того как процесс обучения и оценки модели будет выполнен десять раз (на десяти различных комбинациях тренировочных и тестовых данных), вычисляется средняя эффективность модели по всем блокам.



Крайним случаем CV по k-блокам является метод скользящего контроля — CV по k-блокам с отдельным блоком для каждого примера данных. Таким образом гарантируется использование наибольшего количества данных для обучения модели. Это может показаться полезным, однако данный метод настолько дорог, что редко используется на практике.

Для построения наборов данных для кросс-валидации можно использовать функцию `createFolds()` из пакета `caret`. Как и в случае стратифицированной случайной выборки отложенных данных, функция `createFolds()` попытается сохранить в каждом блоке то же соотношение классов, что и в исходном наборе данных. Ниже приведена команда для создания десяти блоков:

```
> folds <- createFolds(credit$default, k = 10)
```

Результатом функции `createFolds()` является список векторов с номерами строк для каждого из `k=10` блоков, заданных в параметрах

функции. Просмотреть содержимое этих векторов можно с помощью функции `str()`:

```
> str(folds)List of 10$ Fold01: int [1:100] 1 5
12 13 19 21 25 32 36 38 ...$ Fold02: int [1:100]
16 49 78 81 84 93 105 108 128 134 ...$ Fold03:
int [1:100] 15 48 60 67 76 91 102 109 117 123
...$ Fold04: int [1:100] 24 28 59 64 75 85 95 97
99 104 ...$ Fold05: int [1:100] 9 10 23 27 29 34
37 39 53 61 ...$ Fold06: int [1:100] 4 8 41 55 58
103 118 121 144 146 ...$ Fold07: int [1:100] 2 3
7 11 14 33 40 45 51 57 ...$ Fold08: int [1:100]
17 30 35 52 70 107 113 129 133 137 ...$ Fold09:
int [1:100] 6 20 26 31 42 44 46 63 79 101 ...$
Fold10: int [1:100] 18 22 43 50 68 77 80 88 106
111 ...
```

Как видим, первый блок называется `Fold01` и содержит 100 целых чисел, соответствующих 100 строкам из фрейма данных `credit`, которые образуют первый блок. Чтобы создать тренировочные и тестовые наборы данных для построения и оценки модели, необходим еще один шаг. Следующие команды показывают, как создать данные для первого блока. Мы создадим из выбранных 10 % данных тестовый набор и добавим знак «минус», чтобы создать тренировочный набор данных из оставшихся 90 %:

```
> credit01_test <- credit[folds$Fold01, ]>
credit01_train <- credit[-folds$Fold01, ]
```

Для выполнения полной десятиблочной CV нужно повторить этот шаг десять раз, каждый раз сначала создавая модель, а затем вычисляя ее эффективность. В результате будет вычислено среднее значение всех показателей эффективности для получения общей эффективности. К счастью, можно автоматизировать эту задачу, применив несколько методов, рассмотренных ранее.

Чтобы продемонстрировать этот процесс, вычислим каппа-статистику модели дерева решений C5.0 для кредитных данных с использованием десятиблочной CV. Для этого необходимо загрузить следующие R-пакеты: `caret` (для создания блоков), `C50` (для построения дерева решений) и `irr` (для вычисления коэффициента каппа). Последние два пакета выбраны в иллюстративных целях; при желании можно использовать другую модель или другой показатель эффективности, выполняя ту же последовательность операций.

```
> library(caret)> library(C50)> library(irr)
```

Затем мы, как и в прошлый раз, создадим список из десяти блоков. Функция `set.seed()` используется здесь для обеспечения согласованности результатов на случай, если тот же код будет запущен повторно:

```
> set.seed(123) > folds <-  
createFolds(credit$default, k = 10)
```

Наконец, с помощью функции `lapply()` выполним одни и те же операции к списку блоков. Как показано в следующем коде, поскольку не существует готовой функции, которая бы делала то, что нам нужно, мы определим собственную функцию и передадим ее в `lapply()`. Функция разделяет фрейм данных `credit` на тренировочные и тестовые данные, на тренировочных данных строит дерево решений с помощью функции `C5.0()`, генерирует множество прогнозов на основе тестовых данных и сравнивает прогнозируемые и реальные значения, используя функцию `kappa2()`:

```
> cv_results <- lapply(folds, function(x)  
{ credit_train <- credit[-x, ] credit_test  
<- credit[x, ] credit_model <- C5.0(default ~  
., data = credit_train) credit_pred <-  
predict(credit_model,  
credit_test) credit_actual <-  
credit_test$default kappa <-  
kappa2(data.frame(credit_actual,  
credit_pred))$value return(kappa) })
```

Полученные коэффициенты каппа объединяются в список, хранящийся в объекте `cv_results`, который можно просмотреть с помощью `str()`:

```
> str(cv_results) List of 10$ Fold01: num 0.343$  
Fold02: num 0.255$ Fold03: num 0.109$ Fold04: num  
0.107$ Fold05: num 0.338$ Fold06: num 0.474$  
Fold07: num 0.245$ Fold08: num 0.0365$ Fold09:  
num 0.425$ Fold10: num 0.505
```

В процессе десятиблочной CV есть еще один шаг: нужно вычислить среднее этих десяти значений. Несмотря на соблазн ввести `mean(cv_results)`, этого делать не стоит, поскольку `cv_results` не является числовым вектором и результат будет ошибочным. Следует использовать функцию `unlist()`, которая устраняет структуру списка и преобразует `cv_results` в числовой вектор. Теперь мы можем вычислить среднее значение каппа, как и ожидалось:

```
> mean(unlist(cv_results))[1] 0.283796
```

Полученный коэффициент каппа довольно мал, что в соответствии со «справедливой» шкалой интерпретации говорит о том, что данная модель оценки кредитоспособности лишь незначительно лучше, чем случайный выбор. В следующей главе рассмотрим автоматизированные методы, основанные на десятиблочной CV, которые помогают повысить эффективность этой модели.



Пожалуй, многократная k-блочная CV — это на сегодняшний день наилучший метод надежной оценки эффективности модели. Как легко догадаться из названия, этот метод включает в себя многократное применение k-блочной CV и усреднение результатов. Обычная в этом случае стратегия — десятиблочная CV, повторяемая десять раз. Этот метод требует значительных ресурсов, однако дает очень надежную оценку.

Формирование выборок методом начальной загрузки

Менее популярная, но достаточно широко используемая альтернатива k-блочной CV, — это *выборка методом начальной загрузки*, или *бутстреп* (bootstrap sampling). Вообще, так называют статистические методы, в которых для оценки свойств больших наборов данных используются случайные выборки из этих данных. В применении к оценке эффективности модели машинного обучения этот принцип означает создание нескольких случайно выбранных тренировочных и тестовых наборов данных, которые затем используются для оценки статистических показателей эффективности. Затем результаты, полученные для различных случайных наборов данных, усредняются, и получается окончательная оценка будущей эффективности модели.

Чем эта процедура отличается от k-блочной CV? Если при кросс-валидации данные разделяются на блоки, в которых каждый пример может встречаться только один раз, то при методе начальной загрузки в процессе выборки с заменой примеры могут выбираться многократно. Это означает, что из исходного набора данных, состоящего из n примеров, при выполнении бутстрепа будет сформирован один или несколько новых тренировочных наборов данных, каждый из которых также содержит n примеров, некоторые из них повторяются. Затем для каждого из этих тренировочных наборов строится соответствующий тестовый набор данных, в который входят те примеры, которые не вошли в этот тренировочный набор данных.

При описанном использовании выборки с заменой вероятность того, что конкретный пример будет включен в тренировочный набор данных, составляет 63,2 %. Следовательно, вероятность того, что какой-либо пример окажется в тестовом наборе данных, составляет 36,8 %. Другими

словами, тренировочные данные содержат лишь 63,2 % доступных примеров, причем некоторые из них повторяются. В отличие от десятиблочной CV, где для обучения используется 90 % примеров, выборка с заменой менее репрезентативна для полного набора данных.

Поскольку модель, обученная всего лишь на 63,2 % тренировочных данных, вероятно, будет работать хуже, чем модель, обученная на более обширном тренировочном наборе, оценки эффективности метода начальной загрузки могут быть существенно ниже, чем те, что будут получены, когда модель будет обучена на полном наборе данных. Это учитывается в особом варианте метода начальной загрузки, который называется «бутстреп 0,632». В этом случае окончательный показатель эффективности вычисляется как функция эффективности для тренировочных данных (чрезмерно оптимистическая) и тестовых данных (чрезмерно пессимистическая). Затем окончательная частота ошибок вычисляется следующим образом:

$$\text{error} = 0,632 \cdot \text{error}_{\text{test}} + 0,368 \cdot \text{error}_{\text{train}}$$

Одним из преимуществ формирования выборок методом начальной загрузки по сравнению с кросс-валидацией является то, что бутстреп лучше работает с очень маленькими наборами данных. Кроме того, у бутстрепа есть другие области применения, кроме измерения производительности. В частности, из следующей главы мы узнаем, как можно использовать принципы метода начальной загрузки (бутстрепа) для повышения эффективности модели.

Резюме

В этой главе представлен ряд наиболее распространенных показателей и методов оценки эффективности классификационных моделей машинного обучения. Хотя точность является простым средством проверки того, насколько часто модель дает правильные прогнозы, в случае редких событий это может ввести в заблуждение, поскольку реальная частота таких событий может быть обратно пропорциональной частоте их появления в данных.

Показатели, основанные на матрице несоответствий, позволяют лучше сбалансировать штрафы при различных видах ошибок. Внимательное изучение компромиссов между чувствительностью и специфичностью или точностью и полнотой может стать полезным инструментом для исследований последствий ошибок в реальном мире. Для этой цели также полезны визуализации, такие как ROC-кривые.

Стоит также отметить, что иногда лучшим показателем эффективности модели становится анализ того, насколько хорошо она соответствует или не соответствует другим целям. Например, может потребоваться простым языком объяснить логику модели, что исключает некоторые модели из

рассмотрения. Кроме того, даже если модель работает очень хорошо, но слишком медленно или же сложно масштабируется в реальной среде, она будет совершенно бесполезна.

Очевидным дополнением измерения эффективности может стать выявление автоматизированных способов поиска лучших моделей для конкретной задачи. В следующей главе, опираясь на уже проделанную работу, мы изучим способы построения более умных моделей путем систематического, поэтапного уточнения и объединения алгоритмов обучения.

11. Повышение эффективности модели

Когда спортивная команда не справляется со своей задачей — будь то получение золотой олимпийской медали или победа на чемпионате лиги, — ей приходится думать, как улучшить ситуацию. Представьте, что вы тренер такой команды. Как бы вы организовали тренировки? Потребовали бы от спортсменов интенсивнее выкладываться или тренироваться по-другому, чтобы максимально использовать потенциал каждого? Или вы направили бы усилия на то, чтобы улучшить командную работу, более разумно используя сильные стороны спортсменов?

А теперь представьте, что вы тренируете алгоритм на чемпиона мира по машинному обучению. Например, вы рассчитываете принять участие в соревнованиях по интеллектуальному анализу данных, таких как те, результаты которых размещены на сайте Kaggle (<http://www.kaggle.com/>). Или же, возможно, вам просто нужно улучшить работу предприятия. С чего начать? Несмотря на разный контекст, стратегии, которые используются для улучшения результатов спортивной команды, могут применяться и для улучшения показателей статистических алгоритмов обучения.

Ваша задача как тренера заключается в том, чтобы найти такое сочетание методов обучения и навыков командной работы, которое позволило бы достичь поставленных целей. Эта глава основана на материале, рассмотренном в данной книге, и содержит ряд методов повышения эффективности прогнозирования при машинном обучении. Вы узнаете:

- как автоматизировать повышение эффективности модели путем систематического поиска оптимального набора условий обучения;
- какие существуют методы объединения моделей в группы с использованием командной работы для решения сложных задач обучения;
- как применять деревья решений — метод, быстро набравший популярность благодаря своей впечатляющей эффективности.

Ни один из этих методов не будет успешным для всех задач. Тем не менее, изучая записи о победах в соревнованиях по машинному обучению, вы,

вероятно, обнаружите, что для этого был использован хотя бы один из них. Чтобы побеждать в соревнованиях алгоритмов, вам тоже придется освоить эти навыки.

Повышение эффективности готовых моделей

Есть ряд задач обучения, с которыми хорошо справляются стандартные модели, описанные в предыдущих главах. В таких случаях, возможно, нет необходимости тратить много времени на поэтапное уточнение модели; она может сразу работать достаточно хорошо. Но есть и другие задачи, более трудные. Основные понятия могут быть чрезвычайно сложными, требующими понимания некоторых взаимосвязей, или же на задачу могут влиять случайные отклонения, затрудняющие отделение сигнала от шума.

Разработка моделей, которые бы отлично справлялись с подобными задачами, является не только наукой, но в некотором роде и искусством. В одних случаях определить способы повышения эффективности можно с помощью интуиции. В других — методом проб и ошибок. И разумеется, можно прибегнуть к автоматизированным программам.

В главе 5 мы решали сложную задачу: определяли, какие кредиты могут закончиться банкротством. Используя методы настройки эффективности, нам удалось добиться приемлемой точности классификации — около 82 %. Однако из главы 10 вы поняли, что высокая точность иногда вводит в заблуждение. Несмотря на вполне разумную точность, каппа-статистика составила всего примерно 0,28 — это говорит о том, что в действительности модель работает не особенно хорошо. В этом разделе мы вернемся к модели оценки кредитоспособности, чтобы посмотреть, удастся ли улучшить ее результаты.



Для того чтобы выполнить примеры из этой главы, загрузите файл `credit.csv` и сохраните его в рабочем R-каталоге. Загрузите файл в среду R с помощью следующей команды: `credit <- read.csv("credit.csv")`.

Как вы помните, сначала мы использовали для построения классификатора кредитных данных стандартное дерево решений C5.0. Затем повысили его эффективность, настроив параметр `trials`, чтобы увеличить количество повторений бустинга. Увеличив количество итераций от 1 (по умолчанию) до 10, мы смогли повысить точность модели. Такой процесс настройки параметров модели для поиска наилучшего соответствия называется *настройкой параметров*.

Настройка параметров возможна не только для деревьев решений. Например, мы настраивали модели k-NN, когда искали наилучшее значение *k*. Мы также настраивали нейронные сети и метод опорных векторов, где выбирали количество узлов, количество скрытых слоев и

функции ядра. Большинство алгоритмов машинного обучения позволяют регулировать хотя бы один параметр, а самые сложные модели предлагают множество вариантов настройки модели. Это позволяет адаптировать модель к задаче обучения, однако сложность всех возможных вариантов иногда пугает. Нам нужен более систематический подход.

Автоматическая настройка параметров с помощью пакета `caret`

Вместо того чтобы выбирать произвольные значения для каждого из параметров модели — занятия не только утомительного, но и несколько ненаучного, — лучше выполнить поиск среди возможных значений параметров, чтобы найти наилучшую комбинацию.

В пакете `caret`, который мы активно использовали в главе 10, есть инструменты, помогающие автоматизировать настройку параметров модели. Основная функциональность обеспечивается функцией `train()`, которая играет роль стандартизированного интерфейса более чем для 200 различных моделей машинного обучения, применяемых в задачах классификации и регрессии. Используя эту функцию, можно автоматизировать поиск оптимальных моделей, применяя различные методы и показатели для оценки.



Не бойтесь, что моделей окажется слишком много: некоторые из них мы уже рассмотрели в предыдущих главах, а остальные — это лишь простые вариации или расширенные версии базовых концепций. Учитывая полученную информацию, вы сможете понять и остальные методы. Автоматическая настройка параметров подразумевает рассмотрение следующих вопросов.

- Какой тип модели машинного обучения (и какую конкретную реализацию) следует обучать на этих данных?
- Какие параметры модели можно отрегулировать и насколько тщательно их следует настроить для получения оптимальных параметров?
- Какие критерии следует использовать для оценки моделей, чтобы выбрать наилучшую?

Ответ на первый вопрос включает в себя выбор более чем из 200 моделей, доступных в пакете `caret`, такой модели, которая бы лучше всего подходила для решения данной задачи машинного обучения. Очевидно, что это требует широкого и глубокого понимания моделей ML. Эту задачу также можно решить методом исключения.

Почти половину моделей можно исключить сразу, в зависимости от типа задачи: классификация или числовое прогнозирование; многие модели можно исключить из-за формата данных или необходимости избегать

моделей типа «черный ящик» и т.д. В любом случае нет причин, по которым не стоило бы попробовать несколько вариантов, сравнить их результаты и выбрать лучшие.

Ответ на второй вопрос в значительной степени зависит от выбранной модели, поскольку у каждого алгоритма свой, уникальный набор параметров. В табл. 11.1 перечислены параметры настройки для моделей прогнозирования, описанных в этой книге. Помните, что, хотя некоторые из моделей имеют дополнительные опции, не показанные здесь, пакет `caret` поддерживает только те, которые перечислены в табл. 11.1.

Таблица 11.1

Модель	Тип обучающих задач	Имя метода	Параметры
к ближайших соседей	Классификация	knn	k
Наивный байесовский классификатор	Классификация	nb	fL, usekernel
Деревья решений	Классификация	C5.0	model, trials, winnow
Обучение на основе правил OneR	Классификация	OneR	нет
Обучение на основе правил RIPPER	Классификация	JRip	NumOpt
Линейная регрессия	Регрессия	lm	Нет
Регрессионные деревья	Регрессия	rpart	cp
Деревья моделей	Регрессия	M5	pruned, smoothed, rules
Нейронные сети	Двойное назначение	nnet	size, decay
Метод опорных векторов (линейное ядро)	Двойное назначение	svmLinear	C
Метод опорных векторов (радиально-базисное ядро)	Двойное назначение	svmRadial	C, sigma
Случайный лес	Двойное назначение	rf	mtry



Полный список моделей и их параметров настройки, входящих в пакет `caret`, вы найдете в таблице, составленной автором пакета Максом Куном (Max Kuhn), расположенной по адресу <http://topepo.github.io/caret/modelList.html>.

Если вы однажды вдруг забудете параметры настройки для какой-либо модели, можете найти их с помощью функции `modelLookup()`. Для этого достаточно указать имя метода, как показано ниже для модели C5.0:

```
>
modelLookup("C5.0")
  model      parameter
  label forReg forClass probModel1 C5.0
  trials # Boosting
Iterations FALSE      TRUE      TRUE2  C5.0
  model      Model
Type FALSE      TRUE      TRUE3  C5.0      wi
```

`nnow`

`Winnow`

`FALSE`

`TRUE`

`TRUE`

Целью автоматической настройки является поиск набора подходящих моделей в виде матрицы, или *сетки*, сочетаний параметров. Поскольку поиск всех возможных сочетаний нецелесообразен, для построения сетки используется лишь подмножество всех вариантов. По умолчанию `caret` ищет не более трех значений для каждого параметра модели p . Это означает, что будет протестировано не более 3³ подходящих моделей. Например, по умолчанию для алгоритма k-NN в режиме автоматической настройки сравниваются 3¹ = 3 возможные модели со значениями `k=5`, `k=7` и `k=9`. Аналогичным образом настройка дерева решений приведет к сравнению до 27 различных вариантов моделей, включающих в себя сетку из 3³ = 27 сочетаний значений параметров `model`, `trials` и `winnow`. Однако на практике проверяются только 12 моделей. Это связано с тем, что параметры `model` и `winnow` могут принимать только два значения (`tree` или `rules` и `TRUE` или `FALSE` соответственно), так что размер сетки равен $3 \times 2 \times 2 = 12$.



Поскольку сетка поиска, предлагаемая по умолчанию, не всегда подходит для конкретной задачи обучения, `caret` позволяет построить настраиваемую сетку поиска с помощью простой команды, о которой мы расскажем позже. Третий, последний, шаг для автоматической настройки модели — это выбор наилучшей модели из предложенных. Для этого используются методы, рассмотренные в главе 10, включая выбор стратегии повторной выборки для создания тренировочных и тестовых наборов данных, а также использование статистических показателей эффективности модели для измерения точности прогнозирования.

Пакет `caret` поддерживает все стратегии повторной выборки и многие из изученных нами статистических показателей эффективности, в том числе точность и каппа для классификаторов, коэффициент детерминации и RMSE для числовых моделей. При необходимости также можно использовать такие показатели, как чувствительность, специфичность и AUC.

По умолчанию пакет `caret` выбирает из всех вариантов модель с наилучшим значением требуемого показателя эффективности. Поскольку такая методика иногда приводит к выбору моделей, которые обеспечивают лишь небольшой рост эффективности за счет значительного повышения сложности модели, пакет предлагает альтернативные функции выбора модели.

Учитывая большое разнообразие вариантов, хорошо, что многие значения, предлагаемые по умолчанию, выбираются весьма разумно. Например, `caret` будет использовать точность прогнозирования на выборке начальной загрузки (бутстрепа), чтобы выбрать лучший вариант из моделей классификации. Отталкиваясь от этих стандартных значений, можно настроить функцию `train()` для построения большого количества экспериментов.

Простая настройка модели

Чтобы увидеть весь процесс настройки модели, для начала посмотрим, что происходит при настройке модели оценки кредитоспособности с помощью параметров пакета `caret` по умолчанию. Потом зададим параметры модели по своему желанию.

Самый простой способ настройки алгоритма обучения заключается в том, чтобы указать только тип модели с помощью параметра `method`. Поскольку ранее мы уже использовали деревья решений C5.0 для модели оценки кредитов, продолжим работу, оптимизируя этот алгоритм обучения. Основная команда `train()` для настройки дерева решений C5.0 с использованием стандартных значений параметров выглядит следующим образом:

```
> library(caret) > set.seed(300) > m <-  
train(default ~ ., data = credit, method =  
"C5.0")
```

Сначала используется функция `set.seed()` для инициализации генератора случайных чисел R заданным начальным значением. Вероятно, вы помните, мы использовали эту функцию в нескольких предыдущих главах. При выборе начального параметра (в данном случае произвольного числа `300`) генератор будет выдавать случайные числа в предварительно заданной последовательности. Это позволяет повторять симуляции, использующие случайную выборку, с одинаковыми результатами — очень полезная функция, если вы делитесь кодом или пытаетесь повторить предыдущий результат.

Затем с использованием интерфейса R-формулы `default~.` создается дерево. При этом моделируется состояние кредита по умолчанию (`yes` или `no`) с использованием всех остальных признаков из фрейма данных `credit`. Параметр `method="C5.0"` указывает на то, что пакет `caret` будет использовать алгоритм дерева решений C5.0.

Процесс настройки, запущенный предыдущей командой, может занять значительное время (в зависимости от мощности компьютера). Несмотря на сравнительно небольшой набор данных, необходимо выполнить значительный объем вычислений. Среде R нужно будет многократно

генерировать случайные выборки данных, строить деревья решений, вычислять статистические показатели эффективности и оценивать результаты.

Результат эксперимента сохраняется в объекте с именем `m`. Если вы хотите проверить содержимое этого объекта, воспользуйтесь командой `str(m)` — она выведет список всех данных, связанных с объектом. Или же можно просто ввести имя объекта и получить сжатую сводку результатов. Например, в данном случае получим следующее (обратите внимание, что метки в виде цифр добавлены для ясности):

```
1 1000 samples
   16 predictor
   2 classes: 'no', 'yes'

2 No pre-processing
  Resampling: Bootstrapped (25 reps)
  Summary of sample sizes: 1000, 1000, 1000, 1000, 1000, 1000, ...
  Resampling results across tuning parameters:

3 model winnow trials Accuracy Kappa
  rules FALSE 10 0.7147884 0.3181988
  rules FALSE 20 0.7233793 0.3342634
  rules TRUE 10 0.7126357 0.3156326
  rules TRUE 20 0.7225179 0.3342797
  tree FALSE 10 0.7310421 0.3148572
  tree FALSE 20 0.7362375 0.3271043
  tree TRUE 10 0.7285510 0.3093354
  tree TRUE 20 0.7324992 0.3200752

4 Accuracy was used to select the optimal model using the largest value.
  The final values used for the model were trials = 20, model = tree
  and winnow = FALSE.
```

Метки выделяют четыре основных компонента.

1. *Краткое описание входного набора данных.* Если исходные данные вам знакомы и вы правильно применили функцию `train()`, эта информация не должна вызвать удивление.

2. *Отчет о примененных методах предварительной обработки и повторной выборки.* Здесь мы видим, что для обучения моделей было использовано 25 выборок бутстрепа, каждая из которых насчитывала 1000 примеров.

3. *Список рассмотренных предположительно подходящих моделей.* В этом разделе было протестировано 12 моделей, представляющих собой различные сочетания трех параметров настройки C5.0: `model`, `trials` и `winnow`. Для каждой из этих моделей приводятся средняя точность и каппа-статистика.

4. *Выбор наилучшей модели.* Как видно, в итоге была выбрана модель с наибольшей точностью. Это модель C5.0, в которой использовано дерево решений с параметрами `trials=20` и `winnow=FALSE`.

После выбора наилучшей модели функция `train()` использует свои параметры настройки для построения модели на полном наборе входных

данных, который хранится в объекте-списке `m` как `m$finalModel`. В большинстве случаев не приходится напрямую работать с объектом `finalModel`. Достаточно воспользоваться функцией `predict()` для объекта `m`:

```
> p <- predict(m, credit)
```

Полученный вектор прогнозов работает так, как ожидалось. Это позволяет построить матрицу несоответствий, в которой бы сравнивались прогнозируемые и реальные значения:

```
> table(p,
credit)p      no  yes  no  700  2  yes  0
298
```

Из 1000 примеров, использованных для обучения окончательной модели, ошибочно классифицированы всего два. Тем не менее очень важно отметить, что, поскольку модель была построена как на тренировочных, так и на тестовых данных, эта точность является чересчур оптимистичной, поэтому не должна рассматриваться как показатель эффективности для новых данных. Начальная оценка 73 % (показанная в сводных данных модели, сформированных функцией `train()`) является более реалистичной оценкой будущих результатов.

Кроме автоматической настройки параметров, функций `train()` и `predict()`, у пакета `caret` есть еще несколько преимуществ по сравнению с функциями, имеющимися в стандартных пакетах.

Прежде всего, все этапы подготовки данных, выполняемые функцией `train()`, аналогичным образом применяются к данным, используемым для генерации прогнозов. Это включает в себя такие преобразования, как центрирование и масштабирование, а также подстановка отсутствующих значений. Выполняя обработку данных с помощью пакета `caret`, вы гарантируете, что операции, в результате которых была выбрана и оптимизирована наилучшая модель, будут выполнены и при использовании этой модели.

Затем функция `predict()` предоставляет стандартизированный интерфейс для получения прогнозируемых значений класса и прогнозируемых вероятностей классов даже для таких типов моделей, которые обычно требуют дополнительных шагов для получения этой информации. По умолчанию функция работает в режиме прогнозирования классов:

```
> head(predict(m,
credit))[1] no  yes  no  no  yes  noLevels:
no  yes
```


Для того чтобы получить оценку вероятностей для каждого класса, нужно использовать параметр `type="prob"`:

```
> head(predict(m, credit, type =  
"prob" ))  
          no          yes1  0.9606970  0.0  
39302992  0.1388444  0.861155613  1.0000000  0.00  
0000004  0.7720279  0.227972085  0.2948062  0.705  
193856  0.8583715  0.14162851
```

Даже если базовая модель выдает вероятности прогнозирования с помощью другой строки (например, `"raw"` для `naiveBayes`), функция `predict()` автоматически переведет `type="prob"` в соответствующее значение параметра.

Модифицированный процесс настройки

Построенное ранее дерево решений демонстрирует способность пакета `caret` создавать оптимизированную модель при минимальном вмешательстве человека. Параметры настройки, предлагаемые по умолчанию, позволяют легко строить оптимизированные модели. Однако эти стандартные параметры можно изменить на подходящие для конкретной задачи обучения, чтобы выйти на новые, более высокие уровни эффективности.

Дополнительная настройка возможна на каждом этапе выбора модели. Для того чтобы продемонстрировать эту гибкость, изменим наш процесс построения дерева принятия кредитных решений таким образом, чтобы повторить процесс, использованный в главе 10. Как вы помните, мы вычислили каппа-статистику, используя десятиблочную кросс-валидацию (CV). Здесь сделаем то же самое, используя каппа-статистику, чтобы оптимизировать параметр бустинга для дерева решений. Напомню, что бустинг для дерева решений был рассмотрен в главе 5; также мы более подробно о нем поговорим в этой главе.

С помощью функции `trainControl()` создадим набор параметров конфигурации, называемых *управляющим объектом*. Этот объект управляет функцией `train()` и позволяет выбирать критерии оценки эффективности модели, такие как стратегия повторной выборки или показатель, используемый для выбора наилучшей модели. Эту функцию можно использовать для изменения почти всех аспектов эксперимента по настройке, но мы сосредоточимся на двух важных параметрах: `method` и `selectionFunction`.



Для того чтобы получить более подробную информацию, воспользуйтесь командой `?trainControl` и получите список всех параметров.

В функции `trainControl()` параметр `method` используется для выбора метода повторной выборки, такого как выборка отложенных данных или k-блочная CV. В табл. 11.2 приведены допустимые значения параметра `method`, а также все дополнительные параметры для размера выборки и количества итераций. Предлагаемые по умолчанию значения параметров для этих методов повторной выборки соответствуют общепринятым правилам, однако их можно менять в зависимости от размера набора данных и сложности модели.

Таблица 11.2

Метод повторной выборки	Имя метода	Дополнительные параметры и значения по умолчанию
Выборка отложенных данных	LGOCV	$p = 0,75$ (доля тренировочных данных)
k-блочная CV	cv	number = 10 (число блоков)
Многократная k-блочная CV	repeatedcv	number = 10 (число блоков); repeats = 10 (число повторов)
Выборка методом бутстрепа	boot	number = 25 (количество повторных выборок)
Бутстреп 0,632	boot632	number = 25 (количество повторных выборок)
Контроль по отдельным объектам (leave-one-out CV)	LOOCV	Нет

Параметр `selectionFunction` определяет функцию, которая будет выбирать оптимальную модель из предлагаемых кандидатов. Существует три такие функции. Функция `best` просто выбирает вариант с наилучшим значением указанного показателя эффективности. Эта функция используется по умолчанию. Две другие функции применяются для выбора самой экономной или самой простой модели из тех, что находятся в пределах заданного порога максимальной эффективности модели. Функция `oneSE` выбирает простейшую модель, считая, что эффективность модели определяется среднеквадратичной ошибкой, а `tolerance` выбирает простейший вариант в пределах указанного пользователем процента ошибок.



При ранжировании пакетом `caret` моделей по простоте наблюдается некоторая субъективность. О том, как именно выполняется это ранжирование, читайте на странице справки о функциях выбора. Для того чтобы открыть эту страницу, введите `?best` в командной строке R. Для создания управляющего объекта с именем `ctrl` с десятиблочной CV и функцией выбора `oneSE` используйте следующую команду (обратите внимание, что `number=10` здесь указано справочно; поскольку это значение является стандартным для `method="cv"`, его можно опустить):

```
> ctrl <- trainControl(method = "cv", number = 10, selectionFunction = "oneSE")
```

В ближайшее время мы воспользуемся результатом этой функции.

А пока перейдем к следующему шагу нашего эксперимента — построению сетки параметров, которые мы будем оптимизировать. Сетка должна состоять из столбцов со всеми параметрами, настраиваемыми в данной модели, и строк для всех желаемых комбинаций значений параметров. Для дерева решений C5.0 нам понадобятся столбцы с именами `model`, `trials` и `winnow`. Для других моделей машинного обучения обратитесь к таблице, представленной ранее в этой главе, или воспользуйтесь функцией `modelLookup()`, чтобы найти соответствующие параметры, как описано ранее.

Вместо того чтобы заполнять каждую ячейку в этом фрейме данных — что утомительно, если есть много возможных сочетаний значений параметров — воспользуйтесь функцией `expand.grid()`, которая создает фреймы данных из комбинаций всех предоставленных ей значений.

Предположим, что мы хотели бы зафиксировать значения `model="tree"` и `winnow=FALSE` и исследовать восемь разных попыток. Для этого можно создать следующую сетку:

```
> grid <- expand.grid(model = "tree", trials = c(1, 5, 10, 15, 20, 25, 30, 35), winnow = FALSE)
```

Получится фрейм данных сетки, который содержит $1 \times 8 \times 1 = 8$ строк:

```
>
grid  model  trials  winnow1  tree    1  FA
LSE2  tree    5      FALSE3  tree    10  FAL
SE4   tree    15     FALSE5  tree    20  FALS
E6    tree    25     FALSE7  tree    30  FALSE
8     tree    35     FALSE
```

Функция `train()` создает вариант модели для оценки, используя сочетание параметров из каждой строки.

Теперь, имея поисковую сетку и созданный ранее контрольный список, можно приступить к выполнению тщательно настроенного эксперимента с помощью функции `train()`. Как и прежде, зададим генератору случайных чисел начальное значение `300` (выбирается произвольно), чтобы обеспечить повторяемость результатов. Но на этот раз мы передадим функции наш управляющий объект и сетку настройки, добавив параметр `metric="Kappa"`, указывающий статистический параметр,

который будет использоваться функцией оценки модели — в данном случае "oneSE". Полностью команда выглядит следующим образом:

```
> set.seed(300) > m <- train(default ~ ., data =
credit, method = "C5.0", metric =
"Kappa", trControl =
ctrl, tuneGrid = grid)
```

В результате получим объект, который можно просмотреть, введя его имя:

```
1000 samples
 16 predictor
  2 classes: 'no', 'yes'
```

```
No pre-processing
Resampling: Cross-Validated (10 fold)
Summary of sample sizes: 900, 900, 900, 900, 900, 900, ...
Resampling results across tuning parameters:
```

trials	Accuracy	Kappa
1	0.735	0.3243679
5	0.722	0.2941429
10	0.725	0.2954364
15	0.731	0.3141866
20	0.737	0.3245897
25	0.726	0.2972530
30	0.735	0.3233492
35	0.736	0.3193931

```
Tuning parameter 'model' was held constant at a value of tree
Tuning parameter 'winnow' was held constant at a value of FALSE
Kappa was used to select the optimal model using the one SE rule.
The final values used for the model were trials = 1, model = tree
and winnow = FALSE.
```

Большая часть результатов похожа на автоматически настроенную модель, однако есть несколько отличий. Поскольку использовалась десятиблочная CV, то размер выборки для построения каждого варианта модели был сокращен до 900, по сравнению с 1000 примеров, использованных при бутстреппе. Как мы и потребовали, было протестировано восемь вариантов моделей. Кроме того, поскольку параметры `model` и `winnow` были постоянными, их значения не отображаются в результатах — они указаны в сноске.

Теперь лучшая модель весьма существенно отличается от приведенной в предыдущем эксперименте. Если раньше значение `trials` для лучшей модели было равно 20, то теперь `trials=1`. Такое изменение связано с тем, что для выбора оптимальной модели мы использовали правило `oneSE`, а не `best`. По каппа-статистике чистая эффективность модели с `trials=35` выше, однако модель с `trials=1` предлагает почти такую же эффективность при гораздо более простом алгоритме.



Из-за большого количества параметров настройки пакет caret поначалу обескураживает. Пусть это вас не останавливает: не существует более простого способа проверить эффективность моделей, использующих десятиблочную CV. Лучше мысленно разделите эксперимент на две части: объект trainControl(), который определяет критерии тестирования, и сетку настройки, определяющую параметры модели, которые нужно оценить. Передайте все это функции train(), и через некоторое время работы компьютера ваш эксперимент будет завершен!

Повышение эффективности модели с помощью метаобучения

Вместо того чтобы повышать эффективность одной модели, можно объединить несколько моделей в сильную команду. Подобно тому как в лучших спортивных командах навыки игроков не столько совпадают, сколько дополняют друг друга, некоторые лучшие алгоритмы машинного обучения используют группы, составленные из моделей, дополняющих друг друга. Поскольку каждая модель вносит уникальный вклад в задачу обучения, она может хорошо изучить одну группу примеров, но испытывать трудности с другой. Таким образом, разумно используя таланты разных членов команды, можно построить сильную команду из нескольких слабых алгоритмов обучения.

Такая технология объединения нескольких моделей и управления их прогнозами относится к более широкому множеству методов *метаобучения*. Это методы, которые включают в себя обучение тому, как следует учиться. Это может быть что угодно: от простых алгоритмов, которые постепенно повышают эффективность за счет постепенного улучшения структуры (например, автоматической настройки параметров, описанной ранее в этой главе), до очень сложных алгоритмов, в которых используются концепции, заимствованные из эволюционной биологии и генетики, для самоизменения и адаптации к задачам обучения.

Далее в главе мы сосредоточимся именно на метаобучении, поскольку оно касается моделирования взаимосвязей между прогнозами нескольких моделей и желаемым результатом. Рассмотрим достаточно мощные методы командной работы, которые часто используются для создания более эффективных классификаторов.

Понятие ансамблей

Предположим, вы стали участником телевизионной викторины, где разрешается выбрать группу из пяти друзей, которые помогут вам ответить на последний вопрос и получить приз в 1 000 000 долларов. Большинство людей будут пытаться составить группу из экспертов в разных областях. Так, если в эту группу войдут профессора литературы, науки, истории и

искусства, а также эксперт по современной поп-культуре, то выигрыш у вас в кармане. Учитывая их широту знаний, вряд ли вопрос поставит группу в тупик.

В метаобучении используется аналогичный принцип создания команды разнообразных экспертов, которую называют *ансамблем*. В основе всех методов создания ансамблей лежит идея о том, что путем объединения нескольких более слабых алгоритмов обучения можно создать более сильный алгоритм. Методы построения ансамблей различаются главным образом тем, как они отвечают на следующие вопросы.

- Как выбираются и/или строятся слабые модели обучения?
- Как прогнозы, сделанные слабыми алгоритмами, объединяются в единый окончательный прогноз?

Отвечая на них, полезно представить ансамбль в виде диаграммы процесса, как показано на рис. 11.1; почти все ансамблевые подходы следуют такой схеме.

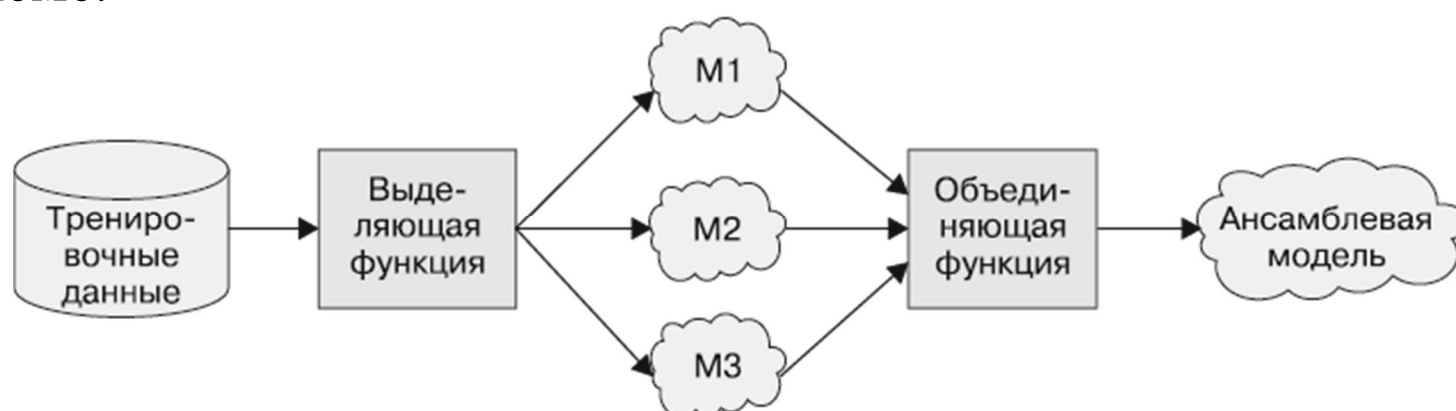


Рис. 11.1. Ансамбли объединяют несколько слабых моделей в одну сильную

Вначале на основе входных тренировочных данных строится ряд моделей. *Выделяющая функция* определяет, сколько тренировочных данных получит каждая модель: полный тренировочный набор данных или лишь выборку, все признаки или только их подмножество.

Идеальный ансамбль включает в себя разнообразный набор моделей, однако выделяющая функция может увеличивать это разнообразие, искусственно изменяя входные данные, чтобы получить различные результаты обучения, даже если модели относятся к одному типу. Например, в ансамбле деревьев решений выделяющая функция может использовать выборку методом бутстрепа для построения уникальных тренировочных наборов данных для каждого дерева или же передавать каждому дереву свои подмножества признаков.

И наоборот, если ансамбль уже включает в себя разнообразный набор алгоритмов, таких как нейронная сеть, дерево решений и классификатор k-NN, то выделяющая функция может передавать каждому алгоритму практически не измененные данные.

После того как модели ансамбля построены, их можно использовать для генерации множества прогнозов, которыми нужно каким-то образом

управлять. *Объединяющая функция* определяет, как разрешать разногласия между прогнозами. Например, ансамбль может определять окончательный прогноз большинством голосов или же использовать более сложную стратегию, такую как взвешивание голосов каждой модели на основе ее предыдущих результатов.

В некоторых ансамблях даже используется другая модель для обучения объединяющей функции на основе различных сочетаний прогнозов. Предположим, что если модели M_1 и M_2 дают ответ «да», то реальным значением класса обычно является «нет». В этом случае ансамбль может научиться игнорировать голоса M_1 и M_2 , если они совпадают. Такой процесс использования прогнозов нескольких моделей для обучения окончательной модели-арбитра называется *стогованием* (рис. 11.2).



Рис. 11.2. Стогование — это сложный ансамбль, в котором используется обучающий алгоритм для объединения прогнозов

Одним из преимуществ ансамблей является то, что они позволяют тратить меньше времени на поиск единственной лучшей модели, так как можно обучить несколько достаточно сильных кандидатов и объединить их. Однако удобство не единственная причина, по которой методы на базе ансамблей продолжают побеждать в соревнованиях по машинному обучению; ансамбли также имеют ряд преимуществ в эффективности по сравнению с отдельными моделями:

- **лучшие способности обобщения для решения будущих задач.** Поскольку единый окончательный прогноз формируется на основе выводов нескольких алгоритмов обучения, ни одно искажение, вносимое отдельными моделями, не может оказаться доминирующим. Это снижает вероятность переобучения;
- **повышенную эффективность для очень больших и очень маленьких наборов данных.** Многие модели сталкиваются с ограничениями по памяти или по сложности, если набор признаков или примеров чрезвычайно большой, поэтому может быть эффективнее обучать несколько моделей, чем одну полную. И наоборот, ансамбли хорошо справляются с очень маленькими наборами данных, поскольку неотъемлемой частью многих ансамблевых структур являются методы многократной выборки, такие как бутстреп. И возможно, самое главное — ансамбль часто можно обучать параллельно, используя методы распределенных вычислений;

- **способность синтезировать данные из разных областей.** Поскольку не существует единого универсального алгоритма обучения, возможность ансамбля использовать способности разных типов алгоритмов обучения становится все более важной, поскольку сложные явления описываются данными, взятыми из разных областей;
- **более тонкое понимание трудных задач обучения.** Явления реального мира часто чрезвычайно сложны, со множеством взаимозависимостей. Модели, которые делят задачу на более мелкие части, очевидно, способны точнее замечать неявные паттерны, которые единая глобальная модель могла бы пропустить.

Ни одно из этих преимуществ не имело бы ценности, если бы невозможно было легко применить ансамблевые методы в R-среде. Для этого есть много пакетов. Рассмотрим самые популярные ансамблевые методы и то, как их можно применять для повышения эффективности разработанной нами модели оценки кредитоспособности.

Бэггинг

Один из первых ансамблевых методов, получивший широкое признание, был основан на технологии *бэггинга* (bagging, сокращенное от bootstrap aggregating). Как писал Лео Брейман (Leo Breiman) в 1994 году, при бэггинге создается несколько тренировочных наборов данных путем выборки методом бутстрепа из исходных тренировочных данных. Затем эти наборы данных используются для генерации моделей на базе единого алгоритма обучения. Прогнозы этих моделей объединяются путем голосования (в случае классификации) или усреднения (в случае числового прогнозирования).



Подробнее о бэггинге читайте в публикации: Breiman L. Bagging predictors. Machine Learning, 1996. Vol. 24. P. 123–140.

Несмотря на то что бэггинг — относительно простой ансамблевый метод, он может работать достаточно хорошо, если состоит из относительно *нестабильных* алгоритмов обучения, то есть таких, которые генерируют модели, имеющие тенденцию существенно изменяться при незначительном изменении входных данных. Нестабильные модели необходимы для обеспечения разнообразия в ансамбле, несмотря на незначительные различия между тренировочными наборами данных, получаемых при бутстрепе. По этой причине при бэггинге часто используются деревья решений, которые имеют тенденцию резко меняться при незначительном изменении входных данных.

Пакет `ipred` представляет собой классическую реализацию бэггинга на основе деревьев решений. Для обучения модели используется функция `bagging()`, аналогичная многим рассмотренным ранее

моделям. Параметр `nbagg` определяет количество деревьев решений в ансамбле, принимающих участие в голосовании (по умолчанию значение этого параметра равно 25). В зависимости от сложности задачи обучения и количества тренировочных данных, увеличение этого числа может повысить эффективность модели до определенного предела. Недостатком модели являются дополнительные вычислительные затраты. Кроме того, обучение большого количества деревьев может занять время.

После установки пакета `ipred` ансамбль создается следующим образом. Мы будем использовать предлагаемое по умолчанию число деревьев решений, равное 25:

```
> library(ipred) > set.seed(300) > mybag <-  
bagging(default ~ ., data = credit, nbagg = 25)
```

Полученная модель, как и ожидалось, совместима с функцией `predict()`:

```
> credit_pred <- predict(mybag, credit) >  
table(credit_pred,  
credit$default) credit_pred no yes no  
699 2 yes 1 298
```

Если сравнить эти результаты с предыдущими, то становится ясно, что эта модель очень хорошо усвоила тренировочные данные. Чтобы узнать, как это повлияет на эффективность алгоритма, можно использовать бэггинг на основе деревьев с десятиблочной CV, применив функцию `train()` из пакета `caret`. Функция, реализующая бэггинг с деревьями в пакете `ipred`, называется `treebag`:

```
> library(caret) > set.seed(300) > ctrl <-  
trainControl(method = "cv", number = 10) >  
train(default ~ ., data = credit, method =  
"treebag", trControl = ctrl) Bagged  
CART1000 samples 16 predictor 2 classes: 'no',  
'yes' No pre-processing Resampling: Cross-Validated  
(10 fold) Summary of sample sizes: 900, 900, 900,  
900, 900, 900, ... Resampling  
results: Accuracy Kappa 0.746 0.3  
540389
```

То, что каппа-статистика для этой модели равна 0,35, говорит о том, что модель бэггинга на основе деревьев работает по меньшей мере так же хорошо, как и лучшее дерево решений C5.0, настроенное нами ранее в этой главе, для которого каппа-статистика была равна 0,32. Это демонстрирует возможности ансамблевых методов: набор простых алгоритмов обучения, работающих совместно, может превзойти весьма сложные модели.

Бустинг

Еще один распространенный метод, основанный на принципе ансамбля, называется *бустингом* (от англ. *boosting* — «форсирование»), потому что он повышает эффективность слабых алгоритмов обучения до показателей, сравнимых с более сильными алгоритмами. Этот метод в значительной степени основан на работах Роберта Шапире (Robert Schapire) и Йоава Фрайнда (Yoav Freund), которые много писали на эту тему.



Подробнее о бустинге читайте в книге: Schapire R.E., Freund Y. *Boosting: Foundations and Algorithms*. Cambridge, MA: The MIT Press, 2012.

Подобно бэггингу, при бустинге используются ансамбли моделей, обученных на многократных выборках данных, с выбором окончательного прогноза путем голосования. Но есть два ключевых различия. Во-первых, повторные выборки наборов данных при бустинге создаются специально таким образом, чтобы генерировать дополняющие друг друга алгоритмы обучения. Во-вторых, не все алгоритмы имеют равные права голоса: при бустинге голос каждого алгоритма обучения имеет вес, основанный на его прошлых результатах. Чем лучше работает модель, тем больше она влияет на окончательный прогноз ансамбля.

Эффективность бустинга зачастую лучше и, во всяком случае, не хуже, чем у наилучших моделей в ансамбле. Поскольку модели в ансамбле подобраны так, чтобы дополнять друг друга, то, чтобы повысить эффективность ансамбля до желаемого уровня, нужно лишь добавить в группу такие классификаторы, чтобы каждый дополнительный классификатор работал лучше, чем в среднем по ансамблю. Поскольку польза этого открытия очевидна, бустинг считается одним из наиболее значимых изобретений в области машинного обучения.



Несмотря на то что бустинг позволяет создать модель, способную обеспечить произвольно низкий уровень ошибок, на практике это не всегда разумно. С одной стороны, по мере ввода дополнительных алгоритмов обучения прирост эффективности уменьшается, так что некоторые уровни эффективности становятся практически недостижимыми. Кроме того, погоня за чистой точностью может привести к тому, что модель окажется перетренированной и не будет делать обобщений для новых данных.

В 1997 году Фрейд и Шапире предложили алгоритм бустинга, получивший название *AdaBoost*, или *адаптивный бустинг*. Этот алгоритм основан на идее генерирования слабых алгоритмов обучения, которые поэтапно изучают большую часть трудноклассифицируемых примеров в

тренировочных данных, уделяя больше внимания (то есть придавая больший вес) тем примерам, которые часто классифицируются неправильно.

Первый классификатор пытается смоделировать результат, начиная с невзвешенного набора данных. Правильно классифицированные примеры с меньшей вероятностью появятся в тренировочном наборе данных для следующего классификатора. И наоборот, трудноклассифицируемые примеры будут появляться чаще.

По мере увеличения количества итераций и добавления новых алгоритмов обучения эти алгоритмы обучаются на наборах данных, последовательно содержащих все более сложные примеры. Процесс продолжается до тех пор, пока не будет достигнута желаемая общая частота ошибок или же пока не повысится эффективность. В этот момент голосу каждого классификатора присваивается вес в соответствии с его точностью на тех тренировочных данных, на которых он был построен.

Принципы бустинга могут применяться практически к любому типу моделей, однако чаще всего они используются для деревьев решений. Мы уже использовали бустинг подобным образом в главе 5 в качестве метода повышения эффективности для дерева решений C5.0.

Еще одна реализация AdaBoost для классификации на основе дерева — это алгоритм *AdaBoost.M1*. Алгоритм *AdaBoost.M1* входит в состав пакета `adabag`.



Подробнее о пакете `adabag` читайте в статье: Alfaro E., Gamez M., Garcia N. `adabag`: An R Package for Classification with Boosting and Bagging // Journal of Statistical Software, 2013. Vol. 54. P. 1–35.

Теперь построим классификатор *AdaBoost.M1* для данных о кредитах. Общий синтаксис этого алгоритма аналогичен другим методам моделирования:

```
> set.seed(300) > m_adaboost <- boosting(default ~ ., data = credit)
```

Как обычно, для прогнозирования к полученному объекту применяется функция `predict()`:

```
> p_adaboost <- predict(m_adaboost, credit)
```

В соответствии с устоявшимся соглашением, вместо того чтобы возвращать вектор прогнозов, эта команда возвращает объект с информацией о модели. Прогнозы хранятся в объекте `class`:

```
> head(p_adaboost$class)[1] "no" "yes" "no"
"no" "yes" "no"
```

В объекте `confusion` содержится матрица несоответствий:

```
>
p_adaboost$confusion
ClassPredicted
Class    no      yes    no      700
         0      yes    0      300
```

Вы заметили, что модель AdaBoost не сделала ни одной ошибки? Но прежде, чем ваши надежды оправдаются, вспомните, что показанная выше матрица несоответствий основана на эффективности модели, вычисленной на основе тренировочных данных. Поскольку бустинг позволяет снизить частоту ошибок до сколь угодно низкого уровня, алгоритм обучения просто продолжал работать, пока не допустил ошибку. Это, скорее всего, привело к переобучению на тренировочном наборе данных.

Для более точной оценки эффективности на неизвестных данных необходимо использовать другой метод оценки. Для десятиблочной CV в пакете `adabag` есть простая функция:

```
> set.seed(300) > adaboost_cv <-
boosting.cv(default ~ ., data = credit)
```

В зависимости от возможностей компьютера выполнение этой функции может занять некоторое время, в течение которого каждая итерация будет выводиться на экран; на моем компьютере последней модели MacBook Pro это заняло около четырех минут. После завершения получается более осмысленная матрица несоответствий:

```
> adaboost_cv$confusion
ClassPredicted
Class    no      yes    no      594      151
         yes    106      149
```

Для того чтобы вычислить каппа-статистику, можно воспользоваться пакетом `vcd`, описанным в главе 10:

```
> library(vcd) >
Кappa(adaboost_cv$confusion)
value
  ASE      z      Pr(>|z|) Unweighted
0.3607 0.0323 11.17 5.914e-
29 Weighted 0.3607 0.0323 11.17 5.914e-29
```

С учетом того, что каппа равна примерно 0,361, это самая эффективная из наших моделей оценки кредитов. Теперь сравним ее с последним ансамблевым методом.



Алгоритм AdaBoost.M1 можно настроить с помощью пакета `caret`, указав параметр `method = "AdaBoost.M1"`.

Случайные леса

Еще один метод, основанный на ансамблях, называемый *случайными лесами* (или *лесами деревьев решений*), использует только ансамбли деревьев решений. Этот метод был предложен Лео Брейманом (Leo Breiman) и Адель Катлер (Adele Cutler). Он сочетает в себе базовые принципы бэггинга со случайным выбором признаков, что позволяет увеличить разнообразие в моделях деревьев решений. После генерации ансамбля деревьев (леса) модель объединяет прогнозы отдельных деревьев путем голосования.



Подробнее о построении случайных лесов читайте в публикации: Breiman L. Random Forests. Machine Learning, 2001. Vol. 45. P. 5–32.

Случайные леса объединяют универсальность и эффективность в единый подход ML. Поскольку ансамбль использует лишь небольшую часть полного набора признаков и эти признаки выбираются произвольно, то случайные леса способны обрабатывать очень большие наборы данных, на которых из-за так называемого «проклятия размерности» применение других моделей привело бы к сбою. В то же время по частоте ошибок для большинства задач обучения случайные леса мало отличаются от других методов.



Термин «случайный лес» зарегистрирован Брейманом и Катлер как товарная марка, однако этот термин иногда используется в разговорной речи для обозначения любого типа ансамбля деревьев решений. Строго говоря, следовало бы использовать более общий термин «лес деревьев решений», за исключением тех случаев, когда речь идет о конкретной реализации Бреймана и Катлер.

Стоит отметить, что, по сравнению с другими ансамблевыми методами, случайные леса весьма конкурентоспособны и имеют ряд важных преимуществ. Например, случайные леса, как правило, проще в использовании и менее подвержены переобучению. В табл. 11.3 перечислены основные преимущества и недостатки моделей случайных лесов.

Таблица 11.3

Преимущества	Недостатки
<p>Универсальная модель, хорошо справляется с большинством задач.</p> <p>Обработывает зашумленные или отсутствующие данные, а также категориальные и непрерывные признаки.</p> <p>Выбирает только самые важные признаки.</p> <p>Может использоваться для данных с чрезвычайно большим количеством признаков или примеров</p>	<p>В отличие от дерева решений эту модель нельзя легко интерпретировать</p>

Благодаря эффективности, универсальности и простоте использования случайный лес — один из самых популярных методов машинного обучения. Позже в этой главе мы детально сравним модель случайного леса с бустингом на основе дерева C5.0.

Обучение случайного леса

В R есть несколько пакетов для создания случайных лесов, однако наиболее точной реализацией алгоритма, разработанного Брейманом и Катлер, является, пожалуй, пакет `randomForest`. Этот пакет также поддерживается `caret`, что обеспечивает возможность автоматической настройки. Синтаксис обучения этой модели выглядит следующим образом.

Синтаксис случайного леса
Использование функции <code>randomForest()</code> из пакета <code>randomForest</code>
<p>Построение классификатора</p> <pre>m <- randomForest(train, class, ntree = 500, mtry = sqrt(p))</pre> <p><code>train</code> – фрейм данных, содержащий тренировочные данные;</p> <p><code>class</code> – факторный вектор, каждый элемент которого содержит класс соответствующей строки тренировочных данных;</p> <p><code>ntree</code> – целое число, определяющее количество деревьев, которые будут сформированы;</p> <p><code>mtry</code> – дополнительный параметр; целое число, определяющее количество признаков, выбираемых случайным образом при каждом разбиении (по умолчанию равно <code>sqrt(p)</code>, где <code>p</code> – число признаков в наборе данных).</p> <p>Функция возвращает объект случайного леса, который затем можно использовать для прогнозирования.</p> <p>Прогнозирование</p> <pre>p <- predict(m, test, type = 'response')</pre> <p><code>m</code> – модель, обученная с помощью функции <code>randomForest()</code>;</p> <p><code>test</code> – фрейм данных, содержащий тестовые данные с тем же набором признаков, что и тренировочные данные, использованные для построения классификатора;</p> <p><code>type</code> принимает одно из трех значений: <code>'response'</code>, <code>'prob'</code> или <code>'votes'</code> и определяет формат вектора прогнозов: прогнозируемый класс, прогнозируемые вероятности или матрица с количеством голосов соответственно.</p> <p>Функция возвращает прогнозы в соответствии со значением параметра <code>type</code>.</p> <p>Пример</p> <pre>credit_model <- randomForest(credit_train, loan_default) credit_prediction <- predict(credit_model, credit_test)</pre>

По умолчанию функция `randomForest()` создает ансамбль из 500 деревьев, которые при каждом разбиении рассматривают `sqrt(p)` случайно выбранных признаков, где `p` — количество признаков в тренировочном наборе данных, а `sqrt()` — R-функция извлечения квадратного корня. Насколько эти параметры, предлагаемые по умолчанию, подходят к конкретной задаче обучения, зависит от характера этой задачи и тренировочных данных. Как правило, более сложные задачи обучения и большие наборы данных (как по количеству признаков, так и по количеству примеров) лучше работают с большим количеством деревьев, хотя необходимо найти баланс с учетом вычислительных затрат на обучение большого количества деревьев.

Цель увеличения количества деревьев — обучить достаточно моделей и предоставить каждому признаку шанс появиться сразу в нескольких моделях. Именно поэтому параметр `mtry` по умолчанию равен `sqrt(p)`; такое значение ограничивает число признаков, чтобы их наборы, выбранные случайным образом, существенно изменялись от дерева к дереву. Например, поскольку данные о кредитах состоят из 16 признаков, каждое дерево в любой момент будет рассматривать только четыре признака.

Посмотрим, насколько предлагаемые по умолчанию параметры `randomForest()` подходят для данных о кредитовании. Мы будем обучать модель так же, как и в случае с другими алгоритмами. Как и раньше, функция `set.seed()` гарантирует воспроизводимость результата:

```
> library(randomForest) > set.seed(300) > rf <-  
randomForest(default ~ ., data = credit)
```

Чтобы просмотреть сводные показатели эффективности модели, достаточно ввести имя полученного объекта:

```
> rfCall:randomForest(formula = default ~ .,  
data = credit)           Type of random  
forest: classification           Number  
of trees: 500 No. of variables tried at each  
split: 4           OOB estimate of error rate:  
23.3% Confusion  
matrix:           no  yes  class.error no  638  62  0.  
08857143 yes 171  129  0.57000000
```

Выходные данные показывают, что, как и ожидалось, случайный лес состоял из 500 деревьев и в каждое разбиение входило по четыре переменные. На первый взгляд, вызывает беспокойство, по-видимому, низкая производительность — судя по матрице несоответствий, частота ошибок составляет 23,3 %, что намного хуже, чем эмпирический риск

любого из рассмотренных нами ранее ансамблевых методов. Однако эта матрица несоответствий не показывает эмпирический риск, а отражает частоту ошибок *out-of-bag* (ООВ) (указанную как `OOBestimateoferrorrate`), которая, в отличие от эмпирического риска, является объективной оценкой ошибки для данного тестового набора. Другими словами, это должна быть достаточно разумная оценка будущих результатов.

Оценка частоты ошибок *out-of-bag* вычисляется во время построения случайного леса. В сущности, для проверки эффективности модели на новых данных можно использовать любой пример, не выбранный для бутстрепа отдельного дерева. После формирования леса любое дерево может делать прогнозы на любом из 1000 примеров в наборе данных, если этот пример не был использован при обучении данного дерева. Затем эти прогнозы подсчитываются для каждого примера и выполняется голосование для определения единого окончательного прогноза для каждого примера. Общая частота ошибок таких прогнозов равна частоте ошибок ООВ.



В главе 10 отмечалось, что вероятность включения любого примера в выборку бутстрепа составляет 63,2 %. Это значит, что при вычислении частоты ошибок ООВ за каждый из 1000 примеров проголосовало в среднем 36,8 % из 500 деревьев случайного леса.

Для того чтобы вычислить каппа-статистику по прогнозам ООВ, можно использовать функцию из пакета `vcd`, как показано далее. В этом коде функция `Kappa()` применяется к первым двум строкам и столбцам объекта `confusion`, в котором хранится матрица несоответствий прогнозов ООВ для объекта `rf` — модели случайного леса:

```
> library(vcd)>
Kappa(rf$confusion[1:2,1:2])          value
  ASE      z  Pr(>|z|)Unweighted
0.381 0.03215 11.85 2.197e-
32Weighted 0.381 0.03215 11.85 2.197e-32
```

Поскольку каппа-статистика равна 0,381, случайный лес — самая эффективная из наших моделей. Это больше, чем у лучшего дерева решений C5.0 с бустингом, для которого каппа составляла примерно 0,325, и больше, чем у модели AdaBoost.M1, для которой каппа равна примерно 0,361. Учитывая этот впечатляющий начальный результат, стоит попытаться оценить его эффективность более формально.

Оценка эффективности случайного леса в условиях моделируемой конкуренции

Как уже отмечалось, функция `randomForest()` поддерживается пакетом `caret`, который позволяет оптимизировать модель, одновременно вычисляя показатели эффективности, помимо частоты ошибок ООВ. Чтобы было еще интереснее, сравним автоматически настроенный случайный лес с самой лучшей автоматически настроенной моделью бустинга на основе C5.0, которую мы построили ранее. Рассмотрим этот эксперимент как попытку выбрать модель для подачи на конкурс по машинному обучению.

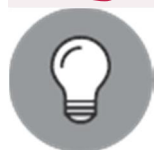
Сначала необходимо загрузить пакет `caret` и выбрать параметры обучения. Для того чтобы сравнение характеристик модели было наиболее точным, мы будем использовать многократную десятиблочную CV, а точнее, десятиблочную CV с десятикратным повторением. Это означает, что на построение моделей уйдет гораздо больше времени, чем ранее, и их оценка будет более затратной с точки зрения вычислений. Однако, поскольку это наше окончательное сравнение, мы **должны** быть уверены, что делаем правильный выбор; победитель этого финального соревнования, выбранный по «лучшему» показателю эффективности, станет единственным кандидатом, которого мы отправим на конкурс по машинному обучению.

Кроме того, мы добавим еще несколько параметров в функцию `trainControl()`. Прежде всего присвоим параметрам `savePredictions` и `classProbs` значение `TRUE`, чтобы сохранить отложенные выборочные прогнозы и прогнозируемые вероятности и впоследствии построить по ним ROC-кривую. Затем присвоим параметру `summaryFunction` значение `twoClassSummary` — это еще одна функция из пакета `caret`, которая вычисляет показатели эффективности, такие как AUC. Полностью описание управляющего объекта выглядит следующим образом:

```
> library(caret) > ctrl <- trainControl(method =  
"repeatedcv", number = 10,  
repeats =  
10, selectionFunction =  
"best", savePredictions =  
TRUE, classProbs =  
TRUE, summaryFunction =  
twoClassSummary)
```

Затем сформируем сетку настройки для случайного леса. Единственный параметр настройки для этой модели — `mtry`. Он определяет, сколько признаков выбирается случайным образом при каждом разбиении. Как мы знаем, по умолчанию случайный лес будет использовать `sqrt(16)` признаков, то есть по четыре признака на каждое дерево. Чтобы повысить точность, протестируем значения вдвое меньше и вдвое больше, а также полный набор из 16 функций. Таким образом, нам нужно построить сетку со значениями 2, 4, 8 и 16:

```
> grid_rf <- expand.grid(mtry = c(2, 4, 8, 16))
```



Случайный лес, который при каждом разбиении рассматривает полный набор признаков, — это, по сути, то же самое, что и бэггинг с моделью дерева решений.

Полученную сетку, а также объект `ctrl` мы передадим в функцию `train()` и используем показатель `"ROC"` для выбора наилучшей модели. Этот показатель соответствует площади под ROC-кривой. Подготовленный эксперимент запускается следующим образом:

```
> set.seed(300) > m_rf <- train(default ~ .,
data = credit, method =
"rf", metric = "ROC", trControl =
ctrl, tuneGrid = grid_rf)
```

Выполнение последней команды может занять некоторое время, поскольку для этого нужно проделать большую работу — на моем MacBook Pro последней модели это заняло около семи минут! Когда обучение случайных лесов закончится, мы сравним наилучший лес с наилучшей моделью бустинга дерева решений, выбранной из деревьев с числом итераций, равным `10`, `25`, `50` и `100`, используя следующий эксперимент из пакета `caret`:

```
> grid_c50 <- expand.grid(model =
"tree", trials = c(10,
25, 50, 100), winnow =
FALSE) > set.seed(300) > m_c50 <- train(default ~
., data = credit, method =
"C5.0", metric = "ROC", trControl
= ctrl, tuneGrid = grid_c50)
```

Когда дерево решений C5.0 наконец завершит работу, можно будет сравнить по пунктам два подхода. Вот результаты для модели случайного леса:

```
> m_rfResampling results across tuning
parameters:mtry ROC Sens Spec 2
0.7579643 0.9900000 0.09766667 4 0.769
5071 0.9377143 0.30166667 8 0.7739714 0
.9064286 0.3863333316 0.7747905 0.8921429
0.44100000
```

А это результаты модели C5.0 с бустингом:

```
> m_c50Resampling results across tuning
parameters:trials ROC Sens Spec
10 0.7399571 0.8555714 0.4346667 25
0.7523238 0.8594286 0.4390000 50 0.7
559857 0.8635714 0.4436667100 0.7566286
0.8630000 0.4450000
```

Как видно из сравнения этих результатов, случайный лес с `mtry=16` определенно является победителем, поскольку его наилучший показатель AUC, равный 0,775, превышает AUC лучшей модели C5.0 с бустингом, равный 0,757.

Чтобы визуализировать эффективность этих методов, можно построить ROC-кривые с помощью пакета `pROC`. Мы передадим функции `roc()` данные о невозвратах кредитов, полученные в ходе тестов (`obs`), а также оценку вероятности значения "yes" для невозвратов кредитов. Обратите внимание: мы используем значения, сохраненные средствами пакета `caret`, и запрашиваем их через функцию `trainControl()`. Затем мы можем использовать функцию `plot()` для построения ROC-кривых:

```
> library(pROC)> roc_rf <- roc(m_rf$pred$obs,
m_rf$pred$yes)> roc_c50 <- roc(m_c50$pred$obs,
m_c50$pred$yes)> plot(roc_rf, col = "red",
legacy.axes = TRUE)> plot(roc_c50, col = "blue",
add = TRUE)
```

Как и ожидалось, полученные кривые показывают, что случайный лес с AUC, равным 0,775, немного превосходит модель C5.0 с бустингом, для которой AUC равен 0,757. На рис. 11.3 видно, что случайному лесу соответствует верхняя кривая.

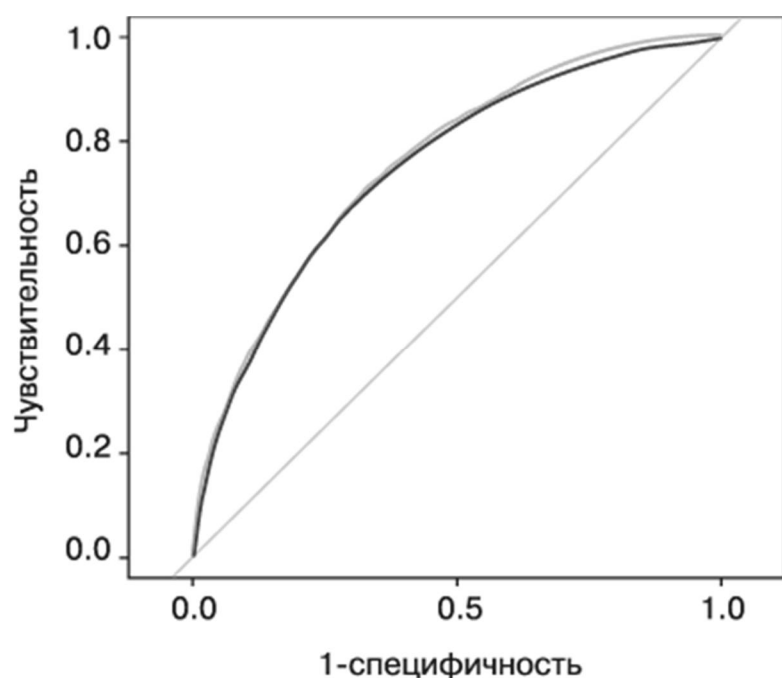


Рис. 11.3. ROC-кривые, позволяющие сравнить случайный лес (верхняя кривая) с бустингом на основе дерева решений C5.0, построенные на наборе данных о невозвратах кредитов

На основании этого эксперимента мы бы отправили на соревнования по машинному обучению именно случайный лес как самую эффективную из наших моделей. Однако до тех пор, пока эта модель не сделает реальные прогнозы на тестовом наборе, выбранном для соревнований, не будет никакого способа узнать точно, победит ли случайный лес в этих соревнованиях. Учитывая наши оценки эффективности, это самый надежный вариант из всех моделей, которые мы оценили, и, если повезет, возможно, мы получим приз.

Резюме

В этой главе вы познакомились с основными методами, которые используются для победы в соревнованиях по интеллектуальному анализу данных и машинному обучению. С одной стороны, автоматические методы настройки помогают выжать максимум эффективности из модели. С другой стороны, повышение эффективности также возможно благодаря созданию совместно работающих групп моделей ML.

Эта глава построена таким образом, чтобы помочь вам подготовить модели для соревнований, однако учтите, что ваши коллеги-конкуренты имеют доступ к тем же самым методам. Не останавливайтесь на достигнутом — продолжайте добавлять собственные методы в список своих приемов. Возможно, вы обладаете уникальными знаниями о предметной области, или, может быть, вашим преимуществом станет внимание к деталям при подготовке данных. В любом случае практика — путь к совершенству, поэтому используйте открытые соревнования, чтобы проверить, оценить и улучшить свои навыки в области машинного обучения.

В следующей, последней, главе мы рассмотрим с высоты птичьего полета способы применения машинного обучения и среды R в некоторых узкоспециализированных и сложных областях. Вы получите знания,

необходимые для использования машинного обучения в решении задач высокой сложности, включающие очень большие или необычные наборы данных.

12. Специальные разделы машинного обучения

Поздравляем с достижением этого этапа в освоении машинного обучения! Если вы еще не начали работу над собственными проектами, то вы сделаете это в ближайшее время. И при этом, возможно, столкнетесь с тем, что преобразовать данные в действия сложнее, чем кажется на первый взгляд.

Собирая данные, вероятно, вы увидите, что информация представлена в нестандартном формате или разбросана по интернет-страницам. Хуже того, потратив много времени на форматирование данных, вы, возможно, обнаружите, что из-за нехватки памяти ваш компьютер начал очень медленно работать. С машиной даже может произойти сбой, или же она может зависнуть. Надеюсь, вас это не остановит, поскольку подобные проблемы поправимы, стоит лишь приложить немного больше усилий.

Эта глава посвящена методам, которые применимы далеко не к любому проекту, однако могут оказаться полезными для решения узкоспециализированных задач. Такая информация будет особенно полезной, если вам часто приходится работать с данными, которые:

- хранятся в неструктурированном виде или в закрытых форматах, таких как веб-страницы, веб-API, электронные таблицы или базы данных;
- относятся к узкоспециализированной области, такой как биоинформатика или анализ социальных сетей;
- слишком велики и не помещаются в памяти или не позволяют выполнить анализ за разумное время.

Не вы одни страдаете от подобных проблем — это проклятие многих исследователей данных, поэтому навыки работы с данными так востребованы. Панацеи от них не существует, однако благодаря совместным усилиям R-сообщества был создан ряд R-пакетов, которые помогают решить некоторые проблемы.

Эта глава — «кулинарная книга» таких решений. Даже если вы ветеран в области R, возможно, вы все равно обнаружите здесь пакет, который упростит ваш рабочий процесс, или же однажды сами создадите пакет, который облегчит работу остальным!

Управление реальными данными и их подготовка

В отличие от примеров, представленных в этой книге, реальные данные редко упаковываются в простой формат CSV, который можно загрузить с веб-сайта. Поэтому приходится приложить немало усилий, чтобы подготовить данные для анализа. Данные необходимо собрать, объединить,

отсортировать, отфильтровать или переформатировать в соответствии с требованиями алгоритма обучения. Этот процесс получил неофициальное название «выпас данных».

Сегодня подготовка данных стала еще важнее, поскольку размер типичных наборов данных вырос с нескольких мегабайтов до нескольких гигабайт, а данные собираются из несвязанных и неупорядоченных источников — зачастую из огромных баз данных. В следующих разделах перечислены некоторые пакеты и ресурсы для поиска и работы с закрытыми форматами и базами данных.

Очистка данных с помощью пакетов tidyverse

В настоящее время формируется новый подход для работы с данными в языке программирования R. Подход быстро занимает доминирующую позицию. Эта новая волна, пропагандируемая Хэдли Уикхемом (Hadley Wickham), разработчиком многих пакетов, вызвавших интерес к R, в настоящее время поддерживается командой RStudio. Настольное приложение RStudio, которое делает язык R существенно более удобным для пользователя, прекрасно интегрируется в новую экосистему, получившую название *tidyverse*, поскольку она предоставляет множество пакетов, предназначенных для упорядочивания (от англ. tidy — «приводить в порядок») данных. Весь набор устанавливается с помощью команды `install.packages("tidyverse")`.

Больше узнать о tidyverse позволяют новые ресурсы, доступные онлайн, начиная с его домашней страницы по адресу <https://www.tidyverse.org>. На странице вы узнаете о различных пакетах, входящих в этот набор, — некоторые из них описаны в данной главе. Кроме того, в Интернете по адресу <https://r4ds.had.co.nz> бесплатно доступна книга Хэдли Уикхема (Hadley Wickham) и Гарретта Гролемунда (Garrett Grolemund) *R for Data Science*, в которой показано, как «безапелляционный» подход tidyverse способен упростить проекты по анализу данных.



Меня часто просят сравнить R и Python с точки зрения науки о данных и машинного обучения. Пожалуй, RStudio и tidyverse являются величайшей ценностью и отличительной чертой R. Возможно, не существует более простых средств, позволяющих анализировать данные.

Обобщение табличных структур данных с помощью tibble

В коллекцию tidyverse входит пакет `tibble` и одноименная структура данных. Слово *tibble* похоже на слово *table* (таблица). Структура данных *tibble* напоминает фрейм данных, однако включает в себя дополнительные функциональные возможности для удобства и простоты использования.

Такие структуры могут применяться почти везде, где применяются фреймы данных. Подробнее о `tibble` читайте в книге *R for Data Science* по адресу <https://r4ds.had.co.nz/tibbles.html>, информация также доступна в среде R, если ввести команду `vignette("tibble")`.

В большинстве случаев использование структур `tibble` является прозрачным и беспроблемным. Однако при необходимости преобразовать `tibble` во фрейм данных используйте функцию `as.data.frame()`. Для обратного преобразования фрейма данных в `tibble` примените функцию `as_tibble()`:

```
> library(tibble)> credit <-  
read.csv("credit.csv")> credit_tbl <-  
as_tibble(credit)
```

Как видно на рис. 12.1, после ввода имени этого объекта вы получите более ясные и информативные сведения, чем для стандартного фрейма данных.

```
> credit_tbl  
# A tibble: 1,000 x 17  
  checking_balance months_loan_dura... credit_history purpose amount savings_balance employment_dura...  
  <fct> <int> <fct> <fct> <int> <fct> <fct>  
1 < 0 DM 6 critical furnitur... 1169 unknown > 7 years  
2 1 - 200 DM 48 good furnitur... 5951 < 100 DM 1 - 4 years  
3 unknown 12 critical education 2096 < 100 DM 4 - 7 years  
4 < 0 DM 42 good furnitur... 7882 < 100 DM 4 - 7 years  
5 < 0 DM 24 poor car 4870 < 100 DM 1 - 4 years  
6 unknown 36 good education 9055 unknown 1 - 4 years  
7 unknown 24 good furnitur... 2835 500 - 1000 DM > 7 years  
8 1 - 200 DM 36 good car 6948 < 100 DM 1 - 4 years  
9 unknown 12 good furnitur... 3059 > 1000 DM 4 - 7 years  
10 1 - 200 DM 30 critical car 5234 < 100 DM unemployed  
# ... with 990 more rows, and 10 more variables: percent_of_income <int>, years_at_residence <int>,  
# age <int>, other_credit <fct>, housing <fct>, existing_loans_count <int>, job <fct>,  
# dependents <int>, phone <fct>, default <fct>
```

Рис. 12.1. Вывод информации об объекте `tibble` более информативен, чем для стандартного фрейма данных

Важно отметить различия между объектами `tibble` и фреймами данных, поскольку при выполнении многих операций `tidyverse` автоматически создает именно объект `tibble`. В целом вы, вероятно, увидите, что `tibble`-объекты ведут себя менее раздражающе, чем фреймы данных. Как правило, `tibble`-объекты делают более разумные предположения о данных, следовательно, вы будете тратить меньше времени на переделку работы после R — например, перекодирование строк в факторы и наоборот. Действительно, главное различие между `tibble`-объектами и фреймами данных заключается в том, что `tibble`-объекты никогда не создаются исходя из предположения, что `stringsAsFactors=TRUE`. Кроме того, `tibble`-объект может также использовать имена столбцов, которые в базовом пакете R недопустимы, такие как ``myvar``, если заключить их в обратные кавычки (```).

Tibble-объекты являются базовым типом объектов в tidyverse и предоставляют дополнительные преимущества для смежных пакетов, описанных в следующих разделах.

Ускорение и упрощение подготовки данных с помощью пакета `dplyr`

Пакет `dplyr` является основой tidyverse. Он предоставляет базовую функциональность, которая дает возможность преобразовывать данные и манипулировать ими. Этот пакет также позволяет начать работу с большими наборами данных в R простым способом. Есть и другие пакеты, обеспечивающие высокую «сырую» скорость или способные обрабатывать еще более массивные наборы данных, однако `dplyr` весьма функционален и является хорошим первым шагом, если нужно выйти за рамки базового R.

В сочетании с tibble-объектами из tidyverse пакет `dplyr` открывает следующие эффектные возможности:

- поскольку этот пакет ориентирован на фреймы данных, а не на векторы, в нем введены новые операторы, позволяющие выполнять обычные преобразования данных с гораздо меньшим количеством кода, при этом сохраняя удобство чтения;
- пакет `dplyr` делает разумные предположения о фреймах данных, которые оптимизируют ваши усилия, а также использование памяти. По возможности этот пакет избегает копирования данных, вместо этого ссылаясь на исходное значение;
- ключевые части кода написаны на C++, что, по мнению авторов пакета, повышает скорость выполнения многих операций в 20–1000 раз по сравнению с базовым R;
- фреймы данных в R ограничены доступной памятью. Пакет `dplyr` позволяет прозрачно связать таблицы с дисковыми базами данных, размер которых может превышать возможности хранения данных в памяти.

Когда вы преодолеете начальную кривую обучения, грамматика `dplyr` станет для вас привычной. В грамматике `dplyr` есть пять ключевых функций, которые выполняют многие из наиболее распространенных преобразований в таблицы данных. Для tibble-объекта доступны следующие преобразования:

- `filter()` — фильтрация строк данных по значениям столбцов;
- `select()` — выборка столбцов данных по имени;
- `mutate()` — преобразование столбцов в новые столбцы путем преобразования значений;
- `summarize()` — агрегирование всех строк данных в одну сводку;
- `arrange()` — упорядочивание строк данных путем сортировки значений.

Эти пять функций `dplyr` можно применять последовательно, используя *конвейерный оператор*. Этот оператор обозначается символами `%>%` и буквально «передает» данные из одной функции в другую. Использование

конвейеров позволяет строить мощные цепочки функций для обработки табличных данных.



Конвейерный оператор является частью пакета `magrittr`, созданного Стефаном Милтоном Бачем (Stefan Milton Bache) и Хэдли Уикхемом (Hadley Wickham). Этот пакет устанавливается по умолчанию вместе с коллекцией `tidyverse`. Название пакета обыгрывает знаменитую картину Рене Магритта «Это не трубка» (она упоминалась в главе 1). Подробнее о пакете `magrittr` читайте на странице `tidyverse` сайта, посвященного этому пакету: <https://magrittr.tidyverse.org>.

Чтобы показать на примере возможности `dplyr`, представьте, что вам нужно обследовать кредитополучателей в возрасте от 21 года, определить среднюю продолжительность кредита (в годах) и сгруппировать результаты по факту возвращения и невозвращения кредита. С учетом этих условий нетрудно понять следующую грамматику `dplyr` — она выглядит почти как пересказ этой задачи на псевдокоде:

```
> credit %>% filter(age >= 21)
%>% mutate(years_loan_duration
= months_loan_duration / 12)
%>% select(default, years_loan_duration)
%>% group_by(default)
%>% summarize(mean_duration =
mean(years_loan_duration))# A tibble: 2 x
2 default mean_duration <fct> <dbl>1
no 1.612 yes 2.09
```

Это лишь один небольшой пример того, насколько последовательности команд `dplyr` способны упростить сложные задачи обработки данных. Это связано с тем, что вследствие более эффективного кода `dplyr` операции зачастую выполняются быстрее, чем эквивалентные команды на базовом R! Полное руководство по `dplyr` выходит за рамки данной книги, однако в Интернете доступно множество учебных ресурсов, в том числе глава из книги *R for Data Science*, размещенная по адресу <https://r4ds.had.co.nz/transform.html>.

Чтение и запись данных во внешние файлы

Что больше всего удручает при анализе данных, так это большой объем работы, необходимой для извлечения и объединения данных из множества закрытых форматов. Огромное количество данных хранится в файлах и базах данных, которые необходимо разблокировать для использования в R. К счастью, есть ряд пакетов, предназначенных именно для этой цели.

Импорт «чистых» таблиц с помощью readr

В состав tidyverse входит пакет `readr` — решение для ускоренной загрузки в R табличных данных, таких как CSV-файлы. Подробнее об этом пакете читайте в книге *R for Data Science*, в главе, посвященной импорту данных в R (<https://r4ds.had.co.nz/data-import.html>), но его основной функционал очень прост.

В пакете доступна функция `read_csv()`, очень похожая на базовую R-функцию `read.csv()`, которая загружает данные из CSV-файлов. Ключевое отличие состоит в том, что функция из tidyverse работает намного быстрее — как утверждают авторы пакета, примерно в десять раз быстрее — и разумнее определяет формат загружаемых столбцов. Например, эта функция обрабатывает числа с символами валюты, преобразовывает столбцы с датами и лучше обрабатывает международные данные.

Для того чтобы создать tibble-объект из CSV-файла, достаточно воспользоваться функцией `read_csv()` следующим образом:

```
> library(readr) > credit <-  
read_csv("credit.csv")
```

При этом используются параметры синтаксического анализа по умолчанию, и они будут отображаться при выводе результатов в R. Эти стандартные значения можно переопределить, передав в `read_csv()` спецификации столбцов, построенные с помощью функции `col()`.

Импорт из файлов Microsoft Excel, SAS, SPSS и Stata с помощью пакета rio

То, что прежде было утомительным и трудоемким процессом, требующим знания конкретных приемов и инструментов из нескольких R-пакетов, сейчас стало элементарным благодаря R-пакету `rio` (сокращение от R input and output — «ввод и вывод в R»). Этот пакет, написанный Чанг-Хонг Чаном (Chung-Hong Chan), Джефффри Ч. Чаном (Geoffrey CH Chan), Томасом Дж. Липером (Thomas J. Leeper) и Кристофером Гандрудом (Christopher Gandrud), называют «швейцарским армейским ножом для ввода/вывода данных». Пакет `rio` позволяет импортировать и экспортировать данные различных файловых форматов, включая (но не ограничиваясь этим) данные, разделенные табуляцией (`.tsv`) и запятыми (`.csv`), JSON, Stata (`.dta`), SPSS (`.sav` и `.por`), Microsoft Excel (`.xls` и `.xlsx`), Weka (`.arff`) и SAS (`.sas7bdat` и `.xpt`).



Полный список типов файлов, которые rio позволяет импортировать и экспортировать, а также более подробные примеры использования этого пакета вы найдете по адресу <http://cran.r-project.org/web/packages/rio/vignettes/rio.html>.

Пакет `rio` состоит из трех функций, предназначенных для работы с закрытыми форматами данных: `import()`, `export()` и `convert()`. Каждая из них делает именно то, что можно ожидать, исходя из ее имени. В соответствии с философией пакета, заключающейся в упрощении работы, каждая функция по расширению имени файла (например, `.csv` или `.xlsx`) определяет тип файла, который требуется импортировать, экспортировать или преобразовать.

Например, чтобы импортировать CSV-файл с данными о кредитах, использовавшийся в предыдущих главах, просто введите:

```
> library(rio) > credit <- import("credit.csv")
```

Как и следовало ожидать, будет создан фрейм данных `credit`. Приятный бонус: нам не только не понадобилось указывать тип CSV-файла, но также `rio` по умолчанию выбирает `stringsAsFactors=FALSE` и другие разумные свойства фрейма данных.

Чтобы экспортировать фрейм данных `credit` в формат Microsoft Excel (`.xlsx`), нужно воспользоваться функцией `export()`, указав желаемое имя файла. Для других форматов просто измените расширение файла на желаемый тип выходных данных:

```
> export(credit, "credit.xlsx")
```

Можно также конвертировать CSV-файл в другой формат напрямую, без импорта, используя функцию `convert()`. Например, следующий код преобразует CSV-файл с кредитными данными в формат Stata (`.dta`):

```
> convert("credit.csv", "credit.dta")
```

Пакет `rio` поддерживает многие распространенные закрытые форматы данных, но не все. В следующем разделе описан другой способ получения данных в R — с помощью запросов к базе данных.

Получение данных путем запросов к базам данных SQL

Большие наборы данных часто хранятся в *системах управления базами данных* (СУБД), таких как Oracle, MySQL, PostgreSQL, Microsoft SQL или SQLite. Эти системы позволяют получать доступ к наборам данных с использованием *языка структурированных запросов* (Structured Query Language, SQL) — языка программирования, предназначенного для извлечения данных из баз данных.

«Чистый» подход к управлению соединениями с базой данных

В RStudio версии 1.1 появился графический способ подключения к базам данных (рис. 12.2). На вкладке соединений (connections) в правой верхней части интерфейса перечислены все существующие соединения с базами данных, найденные в системе. Установка этих соединений обычно выполняется администратором базы данных и зависит от типа базы данных и от операционной системы. Например, для Microsoft Windows, возможно, потребуется установить соответствующие драйверы баз данных, а также использовать приложение ODBC Data Source Administrator; в MacOS и Unix/Linux иногда нужно установить драйверы и отредактировать файл `odbc.ini`.

Полная документация с информацией о возможных типах подключения и инструкциями по установке доступна по адресу <https://db.rstudio.com>.

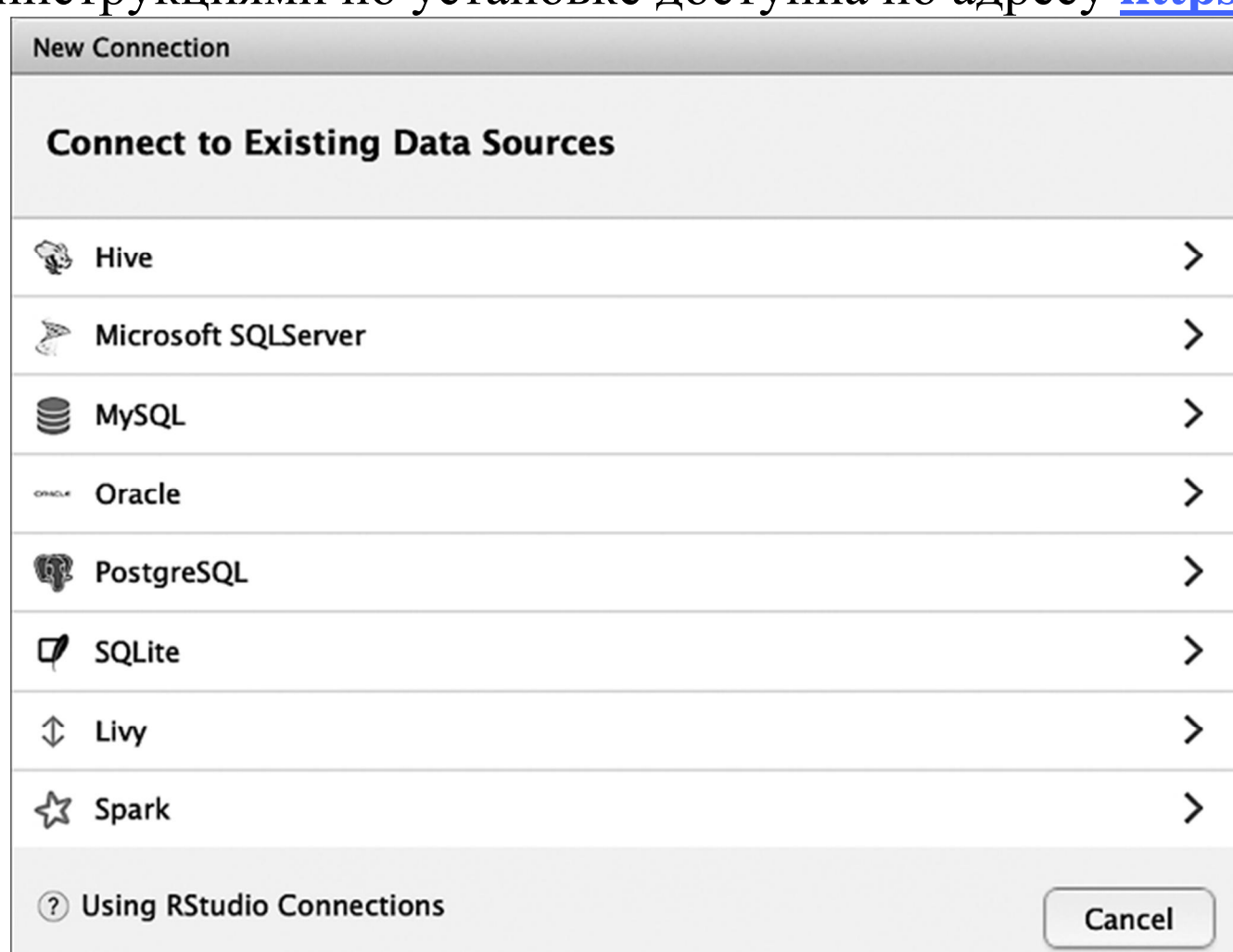


Рис. 12.2. Кнопка New Connection, появившаяся в RStudio v1.1, открывает интерфейс, который помогает подключиться к любым predetermined источникам данных

За фасадом графического интерфейса используется множество R-пакетов для управления соединениями с этими источниками данных. В основе функциональности лежит пакет `DBI`, который предоставляет интерфейс базы данных, совместимый с tidyverse. Пакет `DBI` также управляет драйвером внутренней базы данных, который предоставляется другим R-пакетом. Эти пакеты позволяют подключать R к базам Oracle (`ROracle`), MySQL (`RMySQL`), PostgreSQL (`RpostgreSQL`), SQLite (`RSQLite`) и многим другим.

Чтобы рассмотреть эту функциональность, мы воспользуемся пакетами `DBI` и `RSQLite` для подключения к базе данных SQLite, содержащей уже использованный нами ранее набор данных о кредитах. SQLite — это простая база данных, которая не требует запуска сервера.

Она легко подключается к файлу базы данных, хранящемуся на компьютере, который в данном случае называется `credit.sqlite3`. Перед началом работы не забудьте установить оба требуемых пакета и сохранить файл базы данных в рабочем каталоге R. После этого установите соединение с базой данных с помощью следующей команды:

```
> con <- dbConnect(RSQLite::SQLite(),  
"credit.sqlite3")
```

Чтобы убедиться, что соединение установлено, получите список таблиц базы данных, посмотрите, есть ли среди них таблица `credit`:

```
> dbListTables(con)[1] "credit"
```

С этого момента можно передавать в базу данных команды SQL-запросов и получать в ответ записи в виде R-фреймов данных. Например, чтобы получить всех кредитополучателей в возрасте от 45 лет, нужно направить в базу данных следующий запрос:

```
> res <- dbSendQuery(con, "SELECT * FROM credit  
WHERE age >= 45")
```

Весь набор результатов можно извлечь как фрейм данных с помощью следующей команды:

```
> credit_age45 <- dbFetch(res)
```

Чтобы убедиться, что это сработало, рассмотрим сводную статистику, которая подтверждает, что возраст кредитополучателей начинается с 45 лет:

```
> summary(credit_age45$age)  Min.   1st  
Qu.   Median   Mean   3rd  
Qu.    Max. 45.00   48.00   52.00   53.98   60.00  
75.00
```

После завершения работы желательно очистить набор результатов запроса и разорвать соединение с базой данных, чтобы освободить ресурсы:

```
> dbClearResult(res) > dbDisconnect(con)
```

Кроме SQLite и R-пакетов, предназначенных специально для работы с базами данных, пакет `odbc` позволяет подключать R ко многим другим типам баз данных, используя общий протокол, известный как стандарт *Open Database Connectivity* (ODBC). Стандарт ODBC может применяться независимо от операционной системы или СУБД.

Если вы ранее подключались к базе данных ODBC, то могли ссылаться на нее через *имя источника данных* (Data Source Name, DSN). DSN позволяет устанавливать соединение с базой данных с помощью всего одной строки R-кода:

```
> con <- dbConnect(odbc::odbc(),  
"my_data_source_name")
```

Если настройка более сложная или нужно указать свойства соединения вручную, можно передать функции `dbConnect()` из пакета `DBI` в качестве аргументов полную строку параметров соединения:

```
> library(DBI) > con <-  
dbConnect(odbc::odbc(), 2, databa  
se = "my_database", uid =  
"my_username", pwd =  
"my_password", host =  
"my.server.address", port =  
1234)
```

После установки соединения можно отправлять запросы в базу данных ODBC и получать в ответ таблицы в виде фреймов данных, используя те же функции, которые использовались ранее для примера с SQLite.



Из-за параметров безопасности и брандмауэра инструкции по настройке сетевого подключения ODBC очень зависят от конкретной ситуации. Если возникли проблемы с настройкой параметров соединения, обратитесь к администратору базы данных. Команда RStudio также предоставляет полезную информацию по адресу <https://db.rstudio.com/best-practices/drivers/>.

Доступ к внешней базе данных с помощью пакета `dplyr`

Подключить `dplyr` к внешней базе данных не сложнее, чем использовать его для обычного фрейма данных. Пакет `dbplyr` (`plyr` для баз данных) позволяет использовать любую базу данных, поддерживаемую пакетом `DBI`, в качестве источника данных для `dplyr`. Соединение позволяет извлекать из базы данных `tibble`-объекты. Как правило, достаточно установить пакет `dbplyr`, и тогда `dplyr` сам воспользуется его функциональностью.

Например, подключитесь к использованной нами ранее базе данных SQLite `credit.sqlite3`, а затем сохраните таблицу `credit` в виде `tibble`-объекта с помощью функции `tbl()`:

```
> library(DBI) > con <-  
dbConnect(RSQLite::SQLite(), "credit.sqlite3") >  
credit_tbl <- con %>% tbl("credit")
```

Несмотря на то что пакет `dplyr` перенаправляется в базу данных, объект `credit_tbl` будет работать точно так же, как любой другой `tibble`-объект, и получит все преимущества пакета `dplyr`. Обратите внимание, что, если заменить базу данных SQLite на другую базу,

находящуюся в сети на традиционном SQL-сервере, операции были бы очень похожи.

Например, чтобы получить из базы данных список кредитополучателей в возрасте не менее 45 лет и отобразить сводную статистику по возрасту для этой группы, передайте `tibble`-объект с помощью последовательности функций:

```
> library(dplyr) > credit_tbl %>% filter(age
>= 45) %>% select(age) %>% collect()
%>% summary()      ageMin.      :45.001st
Qu.:48.00Median    :52.00Mean      :53.983rd
Qu.:60.00Max.      :75.00
```

Функции пакета `dbplyr` являются ленивыми — другими словами, функция ничего не делает с базой данных, пока не понадобятся ее результаты. Таким образом, функция `collect()` заставляет `dplyr` получать результаты с сервера в тот момент, когда нужно вычислить итоговые статистические показатели.

При наличии соединения с базой данных многие команды `dplyr` легко преобразуются на сервере в SQL-код. Это означает, что один и тот же R-код можно использовать и в небольших фреймах данных, и для подготовки больших наборов данных, хранящихся в базах данных SQL, — вся тяжелая работа выполняется на удаленном сервере, а не на ноутбуке или настольном компьютере. Таким образом, освоив набор пакетов `tidyverse`, вы сможете применять свой код в любых проектах — от маленьких до огромных.

Традиционный подход к установке SQL-соединения с помощью RODBC

Вместо RStudio и подхода `tidyverse` также можно устанавливать соединение с SQL-сервером с помощью пакета `RODBC`, разработанного Брайаном Рипли (Brian Ripley). Функции пакета `RODBC` получают данные с ODBC-совместимого SQL-сервера и создают в R фрейм данных. Несмотря на то что этот пакет все еще широко используется, он, согласно сравнительным тестам, работает значительно медленнее, чем более новый `odbc`, и приводится здесь в основном для справки.



Описание пакета `RODBC`, которое доступно в R с помощью команды `vignette("RODBC")`, предоставляет обширную информацию о подключении к различным базам данных. Обращайтесь к ней при возникновении проблем.

Чтобы установить соединение `mydb` с базой данных с помощью DSN `my_dsn`, используйте функцию `odbcConnect()`:

```
> library(RODBC) > my_db <-  
odbcConnect("my_dsn")
```

Если же при установке ODBC-соединения требуется имя пользователя и пароль, их нужно указать при вызове функции `odbcConnect()`:

```
> my_db <- odbcConnect("my_dsn", uid =  
"my_username", pwd = "my_password")
```

После того как соединение с базой данных будет установлено, воспользуйтесь функцией `sqlQuery()` и создайте в R фрейм данных из строк базы данных, извлеченных с помощью SQL-запросов.

Функция `sqlQuery()`, как и многие другие функции, создающие фреймы данных, позволяет указать `stringsAsFactors=FALSE`, чтобы предотвратить автоматическое преобразование R символьных данных в факторы.

Функция `sqlQuery()` использует типичные SQL-запросы:

```
> my_query <- "select * from my_table where  
my_value = 1" > results_df <- sqlQuery(channel =  
my_db, query =  
my_query, stringsAsFactor  
s = FALSE)
```

Полученный объект `results_df` представляет собой фрейм данных, который содержит все строки, выбранные SQL-запросом, хранящимся в переменной `my_query`.

Когда закончите использовать базу данных, закройте соединение, как показано в следующей команде:

```
> odbcClose(my_db)
```

Эта команда закрывает соединение `my_db`. Несмотря на то что R автоматически закрывает ODBC-соединения в конце R-сеанса, лучше делать это напрямую.

Работа с онлайн-данными и сервисами

По мере роста объемов данных, доступных из веб-источников, важно обеспечить доступ проектов машинного обучения к интернет-сервисам и взаимодействовать с ними. R с некоторыми оговорками позволяет читать данные из онлайн-источников. Во-первых, по умолчанию R не может получить доступ к защищенным сайтам (использующим протокол <https://> вместо <http://>). Во-вторых, важно отметить, что большинство веб-страниц не предоставляют данные в формате, понятном R. Чтобы такие данные были полезны, необходимо выполнить

Их синтаксический разбор — разбить на части и преобразовать в структурированную форму. Вскоре о способах обойти ограничение будет сказано подробнее.

Однако если ни одно из этих препятствий не имеет места, — другими словами, данные уже размещены в сети на незащищенном сайте и представлены в табличной форме, такой как CSV, которую базовый R позволяет прочитать, то для получения таких данных из Интернета можно использовать R-функции `read.csv()` и `read.table()` — точно так же, как если бы эти данные находились на вашем локальном компьютере. Просто передайте в функцию *полный унифицированный указатель ресурса* (Uniform Resource Locator, URL) набора данных:

```
> mydata <-  
read.csv("http://www.mysite.com/mydata.csv")
```

R также предоставляет функциональные возможности для загрузки из Интернета других файлов, даже если R не может использовать эти файлы напрямую. Для текстового файла используйте функцию `readLines()`:

```
> mytext <-  
readLines("http://www.mysite.com/myfile.txt")
```

Для других типов файлов можно использовать функцию `download.file()`. Чтобы загрузить файл в текущий рабочий каталог R, просто укажите URL-адрес этого файла и папку, где его следует разместить:

```
>  
download.file("http://www.mysite.com/myfile.zip",  
"myfile.zip")
```

Помимо этой базовой функциональности, существует множество пакетов, расширяющих возможности R по работе с онлайн-данными. Самые важные из этих пакетов рассмотрим в следующих разделах. Поскольку сеть огромна и постоянно меняется, это описание далеко от исчерпывающего разбора всех способов подключения R к интернет-источникам данных. Существуют сотни пакетов на все случаи жизни, от узкоспециализированных до массовых проектов.



Наиболее полный и актуальный список пакетов вы найдете на регулярно обновляемом ресурсе CRAN Web Technologies and Services Task View по адресу <http://cran.r-project.org/web/views/WebTechnologies.html>.

Загрузка полного текста веб-страниц

Пакет `RCurl`, разработанный Дунканом Темплом Лангом (Duncan Temple Lang), обеспечивает более надежный способ доступа к веб-страницам,

предоставляя интерфейс R для утилиты curl (клиент для URL-адресов), инструмента командной строки для передачи данных по сети. Программа `curl` — широко распространенный инструмент, работа которого во многом похожа на программируемый веб-браузер; благодаря набору команд эта утилита позволяет получить доступ и загрузить содержимое практически с любого доступного интернет-ресурса. И, в отличие от R, curl может получать доступ к защищенным сайтам, а также передавать инструкции через онлайн-формы. Это невероятно мощная утилита.



Из-за широких возможностей curl полное руководство по использованию этой утилиты выходит за рамки данной главы. Обратитесь к онлайн-документации RCurl по адресу <http://www.omegahat.net/RCurl/>.

Чтобы загрузить страницу после установки и загрузки пакета `RCurl`, достаточно набрать:

```
> packt_page <-  
getURL("https://www.packtpub.com/")
```

Эта команда сохранит полный текст начальной страницы Packt Publishing (включая всю веб-разметку) в символьный R-объект `packt_page`. Как видно в следующем коде, само по себе это не очень полезно:

```
> str(packt_page, nchar.max = 200)chr  
"<!DOCTYPE html>\n <html  
xmlns=\"http://www.w3.org/1999/xhtml\" lang=\"en\"  
xml:lang=\"en\">\n <head>\n <title>Packt  
Publishing| Technology Books, eBooks &  
Videos</title>\" | __truncated__
```

Причина, по которой первые 200 символов страницы выглядят бессмыслицей, заключается в том, что сайты написаны с использованием языка гипертекстовой разметки (Hypertext Markup Language, HTML). Этот язык вставляет в текст страницы специальные теги, которые сообщают веб-браузерам, как отображать данный текст.

Теги `<title>` и `</title>` включают заголовок страницы, таким образом сообщая браузеру, что это домашняя страница Packt Publishing. Аналогичные теги используются для обозначения других частей страницы.

Утилита curl является кросс-платформенным стандартом для доступа к онлайн-контенту, однако, если вы часто работаете с веб-данными в R, лучше подойдет пакет `httr`, разработанный Хэдли Уикхемом. В основе этого пакета лежит `RCurl`, что позволило сделать доступ к веб-данным с HTML-разметкой более удобным и R-совместимым. Вместо того чтобы использовать `RCurl`, внутри пакета `httr` есть собственный

пакет `curl` для получения данных с сайта. Некоторые из различий видны сразу же, стоит лишь загрузить домашнюю страницу Packt Publishing с помощью функции `GET()` из пакета `httr`:

```
> library(httr) > packt_page <-  
GET("https://www.packtpub.com") > str(packt_page,  
max.level = 1) List of 10 $ url : chr  
"https://www.packtpub.com/" $ status_code: int  
200 $ headers : List of 11 ..- attr(*,  
"class")= chr [1:2] "insensitive" "list" $  
all_headers: List of 1 $ cookies : 'data.frame':  
0 obs. of 7 variables: $ content : raw  
[1:162392] 3c 21 44 4f ... $ date :  
POSIXct[1:1], format: "2019-02-24 23:41:59" $  
times : Named num [1:6] 0 0.00372 0.16185  
0.45156... ..- attr(*, "names")= chr [1:6]  
"redirect" "namelookup" "connect" "pretransfer"  
... $ request : List of 7 ..- attr(*,  
"class")= chr "request" $ handle : Class  
'curl_handle' <externalptr>- attr(*, "class")=  
chr "response"
```

Если функция `getURL()` из пакета `RCurl` загружает только HTML, то функция `GET()` из пакета `httr` возвращает, кроме HTML, список со свойствами запроса. Для доступа к самому содержимому страницы нужно использовать функцию `content()`:

```
> str(content(packt_page, type = "text"),  
nchar.max = 200) chr "<!DOCTYPE html>\n<html  
xmlns=\"http://www.w3.org/1999/xhtml\" lang=\"en\"  
xml:lang=\"en\">\n\t<head>\n\t\t<title>Packt  
Publishing | Technology Books, eBooks &  
Videos</title>\n\t\t<script>\n\t\t\tdata"  
____truncated____
```

Чтобы использовать эти данные в R-программе, необходимо обработать данные HTML и привести их в структурированный вид, такой как список или фрейм данных. Необходимые для этого функции будут рассмотрены в следующих разделах.



Подробную документацию и руководство по `httr` вы найдете на веб-странице этого проекта по адресу <https://httr.r-lib.org>. Краткое руководство будет особенно полезно для освоения основных функций.

Синтаксический анализ данных, полученных с веб-страниц

Поскольку HTML-теги многих веб-страниц имеют согласованную структуру, можно создавать программы, которые находят нужные разделы страницы, извлекают их и объединяют в набор данных. Такая практика сбора данных с сайтов и их преобразования в структурированную форму называется *веб-очисткой* (web scraping).



Несмотря на широкое применение, веб-очистку следует рассматривать как последнее средство получения данных из Интернета. Это связано с тем, что любые изменения в базовой структуре HTML могут привести к поломке кода, что потребует усилий для его исправления или, что еще хуже, приведет к внесению в данные незаметных ошибок. Кроме того, в соглашениях об использовании многих веб-сайтов запрещено автоматическое извлечение данных, не говоря уже о том, что трафик программы может перегружать их серверы. Прежде чем приступать к работе над проектом, всегда проверяйте условия использования сайта; вы даже можете обнаружить, что сайт предлагает свои данные бесплатно по соглашению с разработчиком. Вы также можете посмотреть на файл robots.txt — это веб-стандарт, который описывает, какие части сайта могут сканировать боты.

Пакет `rvest` Хэдли Уикхема (в названии пакета обыгрывается слово harvest — «сбор урожая») значительно упрощает очистку веб-страниц, предполагая, что нужные данные находятся в определенном месте HTML-кода.

Начнем с простого примера использования домашней страницы Packt Publishing. Для начала, как и раньше, загрузите страницу, используя функцию `read_html()` из пакета `rvest`. Обратите внимание, что эта функция, если передать ей URL, просто вызывает функцию `GET()` из пакета `httr`:

```
> library(rvest) > packt_page <-  
read_html("https://www.packtpub.com")
```

Предположим, что мы хотим извлечь заголовок страницы; посмотрев на представленный выше HTML-код, мы поняли, что на странице есть только один заголовок, заключенный между тегами `<title>` и `</title>`. Чтобы извлечь этот заголовок, мы передадим имя тега в функцию `html_node()`, а затем еще раз вызовем функцию `html_text()`, которая преобразует результат в простой текст:

```
> html_node(packt_page, "title") %>%  
html_text()[1] "Packt Publishing | Technology  
Books, eBooks & Videos"
```

Обратите внимание на использование конвейерного оператора `%>%`. Как и в базовом пакете `dplyr`, применение конвейеров позволяет строить мощные цепочки функций для обработки данных HTML с помощью функций пакета `rvest`.

Теперь выполним более интересный пример. Предположим, мы хотим выделить список всех пакетов, размещенных на странице CRAN Machine Learning Task View. Как обычно, для начала загрузим HTML-страницу с помощью функции `read_html()`:

```
> library(rvest)> cran_ml <-  
read_html("http://cran.r-  
project.org/web/views/MachineLearning.html")
```

Если просмотреть исходный сайт в браузере, то в одном из разделов обнаружатся интересующие нас данные. Ниже показано только подмножество выходных данных:

```
<h3>CRAN packages:</h3><ul> <li><a  
href=" ../packages/ahaz/index.html ">ahaz</a></li>  
<li><a  
href=" ../packages/arules/index.html ">arules</a></  
li> <li><a  
href=" ../packages/bigrf/index.html ">bigrf</a></li  
> <li><a  
href=" ../packages/bigRR/index.html ">bigRR</a></li  
> <li><a  
href=" ../packages/bmrm/index.html ">bmrm</a></li>  
<li><a  
href=" ../packages/Boruta/index.html ">Boruta</a></  
li> <li><a  
href=" ../packages/bst/index.html ">bst</a></li> <  
li><a  
href=" ../packages/C50/index.html ">C50</a></li> <  
li><a  
href=" ../packages/caret/index.html ">caret</a></li  
>
```

Тегами `<h3>` размечаются заголовки уровня 3, а тегами `` и `` — неупорядоченный список (маркированный, в отличие от упорядоченного, то есть нумерованного, списка) и элементы этого списка соответственно. Требуемые элементы данных заключены в

теги `<a>`, которые являются тегами привязки гиперссылок, указывающих на соответствующую страницу CRAN для каждого пакета.



Поскольку страница CRAN активно обновляется и может измениться в любой момент, не удивляйтесь, если ваши результаты будут отличаться от показанных здесь.

Зная все это, можно извлечь ссылки так же, как и раньше. Единственное исключение состоит в том, что, поскольку мы ожидаем найти несколько результатов, нам нужно вместо функции `html_node()`, которая возвращает только один элемент, использовать функцию `html_nodes()`, чтобы получить вектор результатов. Следующий вызов функции возвращает теги `<a>`, находящиеся внутри тегов ``:

```
> ml_packages <- html_nodes(cran_ml, "li a")
```

Посмотрим на результат с помощью функции `head()`:

```
> head(ml_packages, n = 5){xml_nodeset (5)}[1]
<a href=" ../packages/nnet/index.html">nnet</a>[2]
<a
href=" ../packages/RSNNS/index.html">RSNNS</a>[3]
<a href=" ../packages/rnn/index.html">rnn</a>[4]
<a
href=" ../packages/deepnet/index.html">deepnet</a>
[5] <a
href=" ../packages/RcppDL/index.html">RcppDL</a>
```

Результат включает в себя HTML-теги `<a>` вместе с их содержимым. Чтобы исправить это, достаточно передать результат в функцию `html_text()`. В итоге получим вектор, содержащий имена всех пакетов, перечисленных на странице CRAN Machine Learning Task View, которые затем передаются по конвейеру в функцию `head()`, чтобы отображались только несколько первых значений:

```
> ml_packages %>% html_text() %>%
head()[1] "nnet" "RSNNS" "rnn" "deepnet" "Rc
ppDL" "h2o"
```

Эти простые примеры дают лишь поверхностное представление о возможностях пакета `rvest`. Используя конвейерные операторы, можно искать теги, вложенные в другие теги, или определенные классы HTML-тегов. Чтобы узнать, как решать подобные сложные задачи, обратитесь к документации пакета.



Как правило, очистка веб-страниц — это процесс поэтапного уточнения, поскольку каждый раз приходится определять все более точные критерии, чтобы исключить или, наоборот, включить в результаты все нужные элементы. В самых сложных случаях для достижения абсолютной точности может даже потребоваться участие человека.

Синтаксический анализ XML-документов

XML — это простой текстовый, понятный человеку, но структурированный язык разметки, на котором основано много форматов документов. XML использует структуру тегов, напоминающую HTML, но гораздо сложнее по форматированию. XML является популярным интернет-форматом для хранения структурированных наборов данных.

Пакет `XML`, разработанный Дунканом Темплом Лангом (Duncan Temple Lang), предоставляет набор R-функций, построенный на основе популярного синтаксического анализатора `libxml2`, написанного на языке C и предназначенного для чтения и записи XML-документов. `XML` стал прародителем многих XML-пакетов для синтаксического анализа в среде R, но до сих пор активно используется.



Более подробную информацию о пакете XML, включая простые примеры для быстрого начала работы, вы найдете на сайте проекта по адресу <http://www.omegahat.net/R/XML/>.

В последнее время у библиотеки `libxml2` появился более простой и R-подобный интерфейс — пакет `xml2` от Хэдли Уикхема.

Пакет `xml2` также используется внутри рассмотренного ранее в этой главе пакета `rvest` для синтаксического анализа HTML;

следовательно, `rvest` также можно использовать для анализа XML.



Домашняя страница пакета `xml2` находится по адресу <http://xml2.r-lib.org>. Поскольку синтаксический анализ XML тесно связан с анализом HTML, точный синтаксис здесь не рассматривается. Примеры синтаксического анализа XML вы найдете в документации к этим пакетам.

Синтаксический анализ JSON через веб-API

Онлайн-приложения взаимодействуют между собой с помощью доступных через Интернет функций — *интерфейсов прикладного программирования* (Application Programming Interfaces, API). Эти интерфейсы действуют как обычные веб-

сайты; они получают запрос от клиента по определенному URL и возвращают ответ. Разница в том, что обычный сайт возвращает HTML-код, предназначенный для отображения в браузере, а API обычно возвращает данные, представленные в структурированной форме и предназначенные для машинной обработки.

Найти API на основе XML несложно. Однако, пожалуй, самой распространенной структурой данных для API сегодня является *JavaScript Object Notation (JSON)*. Как и XML, JSON является стандартизированным текстовым форматом, чаще всего используется для структур данных и объектов, передаваемых через Интернет. Формат стал популярным в последнее время благодаря своему происхождению из браузерных приложений, написанных на JavaScript. Однако, несмотря на такую «родословную», его полезность не ограничивается Интернетом. Простота, с которой структуры данных JSON могут быть поняты людьми и проанализированы машинами, делает эти структуры данных привлекательными для самых разных проектов.

JSON основан на формате типа `{ключ: значение}`. Фигурные скобки `{ }` заключают в себе JSON-объект, а ключ и значение обозначают свойство объекта и статус этого свойства. Объект может иметь любое количество свойств, а сами свойства могут, в свою очередь, быть объектами. Например, объект JSON для этой книги может выглядеть примерно так:

```
{ "title": "Machine Learning with  
R", "author": "Brett Lantz", "publisher":  
{ "name": "Packt Publishing", "url":  
"https://www.packtpub.com" }, "topics": ["R",  
"machine learning", "data mining"], "MSRP":  
54.99 }
```

Пример иллюстрирует типы данных, доступные в JSON: числовые, символьные, массивы (заключенные между символами `[` и `]`) и объекты. Здесь не показаны значения `null` и `Boolean` (`true` или `false`). Передача этих типов объектов из одного приложения в другое и из приложения в браузер обеспечивает работу многих популярных сайтов.



Подробнее о формате JSON читайте на сайте <http://www.json.org/>.

Существует ряд пакетов, которые позволяют преобразовывать данные в формат JSON и обратно. Пакет `jsonlite`, созданный Йеруном Омсом (Jeroen Ooms), быстро стал популярным, поскольку он позволяет создавать структуры данных, более согласованные и свойственные R, чем другие пакеты, особенно при использовании данных из веб-API. Подробнее о том,

как использовать этот пакет, читайте на его странице GitHub по адресу <https://github.com/jeroen/jsonlite>.

После установки пакета `jsonlite` для преобразования R-объекта в строку JSON нужно использовать функцию `toJSON()`. Обратите внимание, что в результатах этой функции символы кавычек экранированы с помощью символов `\`:

```
> library(jsonlite) > ml_book <- list(book_title = "Machine Learning with R", author = "Brett Lantz") > toJSON(ml_book) {"book_title":["Machine Learning with R"],"author":["Brett Lantz"]}
```

Чтобы преобразовать строку JSON в R-объект, используйте функцию `fromJSON()`. Кавычки в строке должны быть экранированы:

```
> ml_book_json <- "{ \"title\": \"Machine Learning with R\", \"author\": \"Brett Lantz\", \"publisher\": { \"name\": \"Packt Publishing\", \"url\": \"https://www.packtpub.com\" }, \"topics\": [\"R\", \"machine learning\", \"data mining\"], \"MSRP\": 54.99}" > ml_book_r <- fromJSON(ml_book_json)
```

В результате получим список в форме, очень похожей на JSON:

```
> str(ml_book_r) List of 5 $ title : chr "Machine Learning with R" $ author : chr "Brett Lantz" $ publisher: List of 2 .. $ name: chr "Packt Publishing" .. $ url : chr "https://www.packtpub.com" $ topics : chr [1:3] "R" "machine learning" "data mining" $ MSRP : num 55
```



Подробнее о пакете `jsonlite` читайте в публикации: Ooms J. The `jsonlite` Package: A Practical and Consistent Mapping Between JSON Data and R Objects, 2014, доступной по адресу <http://arxiv.org/abs/1403.2805>.

Открытые API позволяют программам, подобным R, систематически обращаться к сайтам и получать результаты в формате JSON с использованием таких пакетов, как `Rcurl` и `httr`. Почти все веб-сайты, которые содержат интересные данные, предлагают API для запросов, хотя некоторые из них взимают за это плату и требуют ключ доступа. Среди популярных примеров — API для Twitter, Google Maps и Facebook. Полное

руководство по использованию веб-API составило бы отдельную книгу, однако сам процесс состоит всего из нескольких шагов — сложности начинаются при углублении в детали.

Предположим, что мы хотим запросить Apple iTunes API, чтобы найти альбомы The Beatles. Для этого нам сначала нужно просмотреть документацию по iTunes API, которая находится по адресу <https://affiliate.itunes.apple.com/resources/documentation/itunes-store-web-service-search-api/>, чтобы сформировать URL-адрес и параметры, необходимые для выполнения этого запроса. Затем мы передадим эту информацию в функцию `GET()` пакета `httr`, добавив список параметров запроса для отправки на URL поиска:

```
> library(httr) > music_search <-  
GET("https://itunes.apple.com/search",  
    query = list(term =  
"Beatles",  
media =  
"music",  
entity =  
"album",  
limit = 10))
```

После ввода имени полученного объекта мы увидим некоторые детали запроса:

```
> music_searchResponse  
[https://itunes.apple.com/search?term=Beatles&media=music&entity=album&limit=10] Date: 2019-02-25  
00:33 Status: 200 Content-Type:  
text/javascript; charset=utf-8 Size: 9.75  
kB { "resultCount": 10, "results":  
[ { "wrapperType": "collection",  
"collectionType": "Album",  
"artistId": 136975, "collectionId": 402060584,  
"amgArtistId": 3644, "artistName": "The B...
```

Чтобы получить доступ к JSON, воспользуемся функцией `content()`. Затем результат преобразуем в R-объект с помощью функции `fromJSON()` из пакета `jsonlite`:

```
> library(jsonlite) > music_results <-  
fromJSON(content(music_search))
```

Объект `music_results` представляет собой список, содержащий данные, возвращенные из API iTunes. Результатов слишком много, чтобы приводить их здесь целиком, однако

команда `str(music_results)` позволяет увидеть структуру этого объекта. Как видим, эти интересные данные хранятся внутри объекта `results`. Например, вектор названий альбомов находится по имени `collectionName`:

```
> music_results$results$collectionName[1] "The Beatles Box Set"[2] "Abbey Road"[3] "The Beatles (White Album)"[4] "The Beatles 1967-1970 (The Blue Album)"[5] "1 (2015 Version)"[6] "Sgt. Pepper's Lonely Hearts Club Band"[7] "The Beatles 1962-1966 (The Red Album)"[8] "Revolver"[9] "Rubber Soul"[10] "Love"
```

Затем эти элементы данных можно использовать в R-программе любым желаемым способом.



Поскольку Apple iTunes API обновляется, вы можете столкнуться с тем, что ваши результаты отличаются от представленных здесь.

Работа со специфическими данными

Машинное обучение, что естественно, применяется для решения задач в самых разных областях. Его основные методы везде одинаковы, однако некоторые из них настолько специализированы, что сформировались особые сообщества для разработки решений задач, свойственных отдельным предметным областям. Это приводит к открытию новых методов и появлению терминологии, которая имеет отношение только к специфическим задачам, свойственным данной предметной области.

В этом разделе рассмотрено несколько областей, в которых широко используются методы ML, но для того, чтобы полностью раскрыть их потенциал, необходимы специальные знания. Поскольку на подобные темы написаны целые книги, здесь приведен лишь краткий обзор. Для получения более подробной информации обратитесь к ресурсам, указанным в каждом подразделе.

Анализ данных в биоинформатике

Биоинформатика — это наука, где компьютеры и анализ данных применяются для решения биологических проблем, особенно для лучшего понимания генома. Поскольку генетические данные резко отличаются от большинства других типов данных, анализ данных в области биоинформатики ставит ряд уникальных задач. Например, поскольку живые существа имеют огромное количество генов, а генетическое секвенирование все еще стоит относительно дорого, типичные наборы

генетических данных в ширину гораздо больше, чем в длину; другими словами, в них признаков (генов) больше, чем примеров (организмов, которые были секвенированы). Это создает проблемы при попытке применить к таким данным обычные методы визуализации, статистические тесты и методы машинного обучения. Кроме того, использование патентованных методик типа «лаборатория на чипе» с *микромассивами* требует узкоспециализированных знаний уже на этапе простой загрузки генетических данных.



На сайте CRAN представлены список задач и некоторые специализированные R-пакеты для статистической генетики и биоинформатики: <http://cran.r-project.org/web/views/Genetics.html>.

Решением некоторых подобных проблем занимается проект *Bioconductor* Центра исследований рака Фреда Хатчинсона в Сиэтле, штат Вашингтон. Целью проекта является предоставление стандартизированного набора методов анализа геномных данных. Используя R как основу, Bioconductor добавляет пакеты, относящиеся к биоинформатике, и документацию, относящуюся к базовому программному обеспечению R.

Bioconductor предоставляет рабочие процессы для анализа ДНК и данных белковых микрочипов, получаемых с таких распространенных платформ, как Affymetrix, Illumina, NimbleGen и Agilent. Дополнительная функциональность включает в себя аннотирование последовательностей, процедуры тестирования, специализированные визуализации, учебные пособия, документацию и многое другое.



Подробнее о проекте Bioconductor читайте на сайте проекта по адресу <http://www.bioconductor.org>.

Анализ и визуализация сетевых данных

Данные, получаемые из социальных сетей, и наборы графовых данных состоят из структур, которые описывают зависимости, или *связи* (иногда называемые *ребрами*), между людьми или объектами — *узлами*. При наличии N узлов можно построить матрицу потенциальных связей размера $N \times N = N^2$. Увеличение количества узлов значительно усложняет вычислительный процесс.

Сетевой анализ занимается статистическими показателями и способами визуализации, которые позволяют идентифицировать значимые связи. Например, на рис. 12.3 показаны три кластера узлов, обозначенные на рисунке кругами, которые соединены через квадратный узел,

расположенный в центре. Сетевой анализ позволяет выявить важность квадратного узла по отношению к другим ключевым показателям.

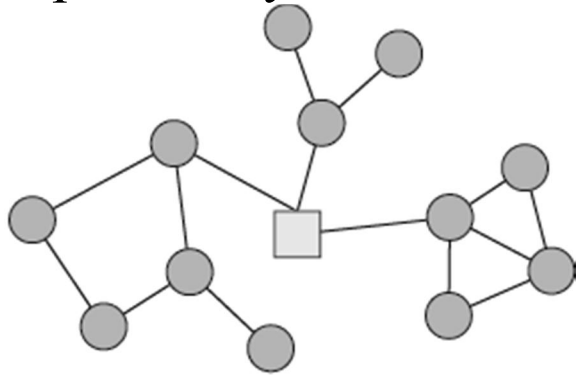


Рис. 12.3. Социальная сеть, в которой узлы сгруппированы в три кластера, расположенных вокруг центрального узла (обозначенного квадратом)

Пакет `network`, разработанный Картером Т. Баттсом (Carter T. Butts), Дэвидом Хантером (David Hunter) и Марком С. Хэндкоком (Mark S. Handcock), предлагает специализированную структуру данных для работы с сетями. Эта структура данных необходима, так как матрица для хранения N^2 потенциальных связей быстро исчерпывает память; в структуре данных `network` используется разреженное представление для хранения только существующих связей. Это значительно экономит память, если большинство потенциальных зависимостей в действительности не существует. Тесно связанный с `network` пакет `sna` (от Social Network Analysis — анализ социальных сетей) позволяет анализировать и визуализировать объекты `network`.



Более подробно о пакетах `network` и `sna`, включая руководства и документацию, читайте на сайте проекта, размещенного Вашингтонским университетом по адресу <http://www.statnet.org/>. Отличные учебные материалы также предлагает лаборатория анализа социальных сетей из Стэнфордского университета по адресу <https://sna.stanford.edu/rlabs.php>. Еще один набор инструментов для визуализации и анализа сетевых данных представлен в пакете `igraph` Габора Чарди (Gabor Csardi). Этот пакет позволяет обрабатывать и вычислять статистические показатели для очень больших сетей. Преимуществом пакета `igraph` является наличие аналогичных пакетов для языков программирования Python и C, что позволяет использовать этот пакет практически везде, где выполняется анализ. Как вы вскоре убедитесь, этот пакет очень прост в использовании.



Подробнее о пакете `igraph`, включая демонстрационные примеры и учебные пособия, читайте на его домашней странице по адресу <http://igraph.org/r/>.

Анализ сетевых данных в R требует использования специализированных форматов, поскольку сетевые данные обычно не хранятся в типичных табличных структурах, таких как CSV-файлы и фреймы данных. Как уже отмечалось, поскольку между n сетевыми узлами существует n^2 потенциальных связей, то табличная структура быстро станет громоздкой для любых значений n , кроме самых маленьких. Данные графов хранятся в формате, в котором учитываются только действительно существующие связи; о том, каких связей нет, можно догадаться по отсутствию данных о них.

Возможно, самый простой из таких форматов — это *список ребер*, представляющий собой текстовый файл, в котором каждая строка соответствует существующей в сети связи. Каждому узлу присваивается уникальный идентификатор, а связь между узлами описывается как строка, в которой размещены идентификаторы соединенных узлов через пробел. Например, следующий список ребер описывает три связи узла 0 с узлами 1, 2 и 3.

```
0 10 20 3
```

Для того чтобы загрузить сетевые данные в R-пакет `igraph`, нужно воспользоваться функцией `read.graph()`. Она читает файлы со списками ребер, а также другие, более сложные форматы, такие как *Graph Modeling Language* (GML). Чтобы наглядно показать эту функциональность, воспользуемся набором данных, описывающим дружеские отношения между членами небольшого клуба карате. Для выполнения этого примера загрузите файл `karate.txt` с веб-сайта Packt Publishing и сохраните его в рабочем каталоге R. Затем установите пакет `igraph`, после чего сеть клуба карате в R можно прочитать следующим образом:

```
> library(igraph)> karate <-  
read.graph("karate.txt", "edgelist", directed =  
FALSE)
```

Будет создан объект разреженной матрицы, который можно использовать для построения графов и сетевого анализа. Обратите внимание, что параметр `directed=FALSE` создает сеть с ненаправленными или двунаправленными связями между узлами.

Поскольку набор данных о клубе карате описывает дружеские отношения, это означает, что если `Person1` дружит с `Person2`, то `Person2` также дружит с `Person1`. И наоборот, если бы в наборе данных были описаны результаты боев, то факт, что `Person1` победил `Person2`, конечно же, не означает, что `Person2` победил `Person1`. В этом случае должен быть установлен параметр `directed=TRUE`.



Используемый здесь набор данных о сети клуба карате был составлен М.И. Дж. Ньюманом (М. E.J. Newman) из Мичиганского университета. Впервые этот набор данных был представлен в статье: Zachary W.W. An Information Flow Model for Conflict and Fission in Small Groups // Journal of Anthropological Research, 1977. Vol. 33. P. 452–473.

Чтобы изучить этот граф, воспользуемся функцией `plot()`:

```
> plot(karate)
```

Что получим в результате, показано на рис. 12.4.

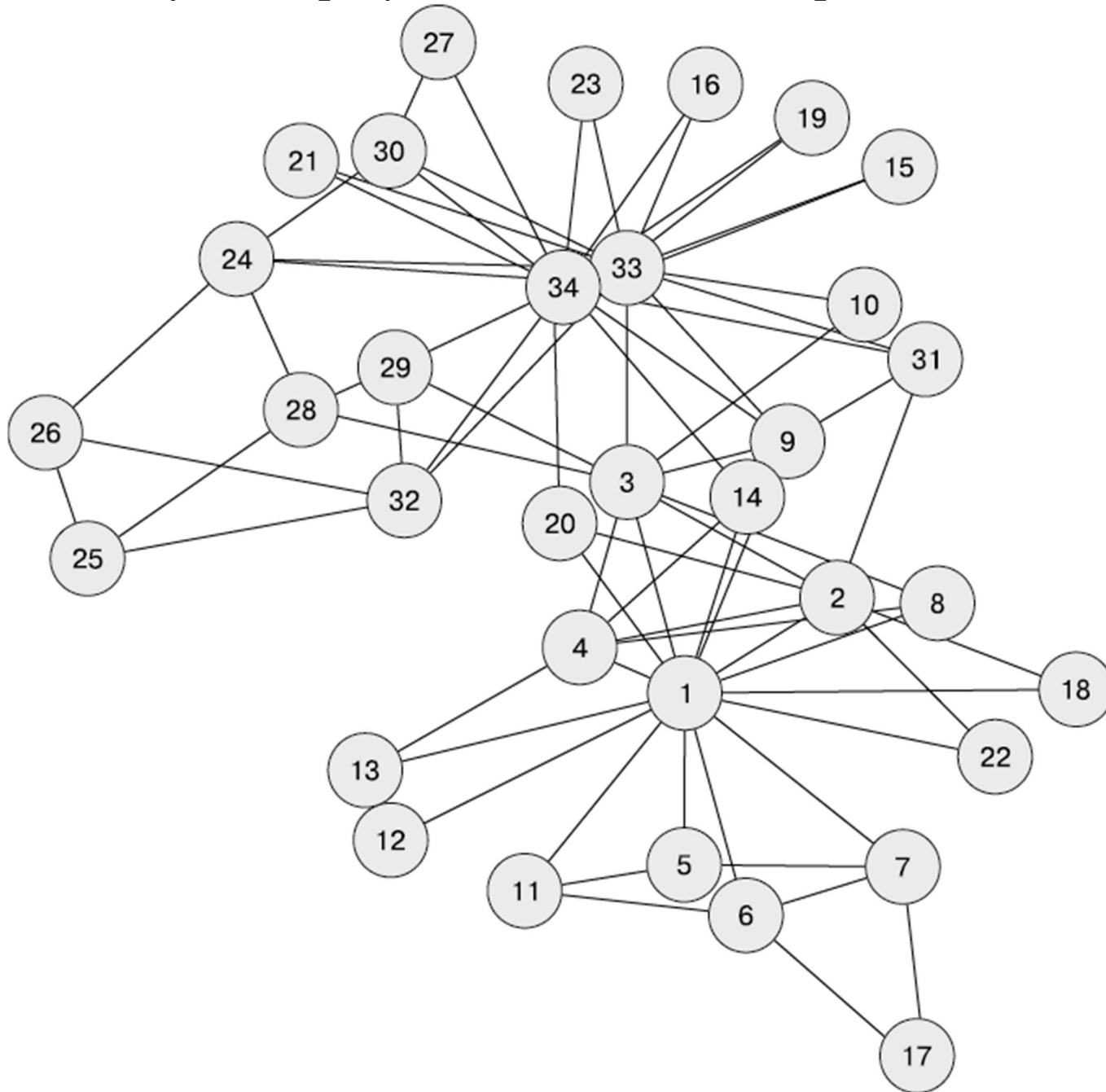


Рис. 12.4. Сетевая визуализация набора данных о клубе карате. Связи указывают на дружеские отношения

Как видно из этой сетевой визуализации, в данном карате-клубе есть несколько членов, имеющих большое количество связей с другими членами клуба. Узлы 1, 33 и 34 расположены ближе к центру, чем другие узлы, которые находятся на периферии графа клуба.

Используя `igraph` для вычисления показателей графа, можно подтвердить нашу догадку аналитически. *Степень* узла означает количество узлов, с которыми он связан.

Функция `degree()` подтверждает наше предположение о том, что узлы 1, 33 и 34 имеют больше связей, чем другие узлы сети, имея 16, 12 и 17 связей соответственно:

```
> degree(karate)[1] 16 9 10 6 3 4 4 4 5 2
3 1 2 5 2 2 2 2[19] 2 3 2 2 2
5 3 3 2 4 3 4 4 6 12 17
```

Поскольку некоторые соединения более важны, чем другие, для оценки связности узлов разработано множество сетевых показателей. Сетевой показатель, называемый *центральностью по посредничеству*, определяет количество кратчайших путей между узлами, которые проходят через каждый узел. Узлы, которые действительно являются наиболее важными для всего графа, имеют более высокое значение центральности, поскольку они служат мостами между другими узлами. С помощью функции `betweenness()` можно получить вектор показателей центральности:

```
> betweenness(karate)[1]
231.0714286 28.4785714 75.8507937 6.2880952[5]
0.3333333 15.8333333 15.8333333 0.0000000[9]
] 29.5293651 0.4476190 0.3333333 0.0000000[13]
0.0000000 24.2158730 0.0000000 0.0000000
0[17] 0.0000000 0.0000000 0.0000000
17.1468254[21] 0.0000000 0.0000000 0.0000000
0 9.3000000[25] 1.1666667 2.0277778 0.0000000
11.7920635[29] 0.9476190 1.5428571 7.6095238
8 73.0095238[33] 76.6904762 160.5515873
```

Поскольку узлы 1 и 34 имеют большие значения центральности, чем другие узлы, они являются более важными для сети дружбы карате-клуба. Эти два человека с обширными личными дружескими связями могут быть тем «клеем», на котором держится весь клуб.



Центральность по посредничеству — всего один из многих показателей, предназначенных для определения важности узла, и даже не единственная мера центральности. Подробнее о других свойствах сетей читайте в документации `igraph`.

Пакеты `sna` и `igraph` позволяют вычислять множество показателей графов, которые затем можно использовать как исходные данные для функций машинного обучения. Предположим, что мы хотим построить модель, предсказывающую, кто победит на выборах президента клуба. Тот

факт, что узлы 1 и 34 имеют много связей, говорит о том, что они могут иметь социальный капитал, достаточный для руководящей роли. Данные о связях могут быть очень ценными для предсказания результатов выборов.



Комбинируя анализ сети с машинным обучением, такие сервисы, как Facebook, Twitter и LinkedIn, предоставляют обширные хранилища сетевых данных для прогнозирования поведения пользователей. Ярким примером этого является президентская кампания 2012 года в США, в ходе которой главный специалист по данным Рейд Гани использовал данные Facebook, чтобы идентифицировать людей, которых можно убедить проголосовать за Барака Обаму.

Повышение эффективности R

Базовый R имеет репутацию (в некоторой степени заслуженную) медленного языка с неэффективным использованием памяти. На современных компьютерах и для наборов данных, насчитывающих несколько тысяч записей, эти недостатки, как правило, незаметны. Однако наборы данных из миллиона и более записей могут превышать пределы возможностей компьютерного оборудования потребительского уровня. Проблема становится еще серьезнее, если набор данных содержит много признаков или если используются сложные алгоритмы обучения.



В CRAN есть описания задач, подразумевающих высокопроизводительные вычисления, и перечисляются пакеты, расширяющие пределы возможностей R: <http://cran.r-project.org/web/views/HighPerformanceComputing.html>.

Сейчас быстро развиваются пакеты, которые расширяют возможности R сверх базовых. Эта работа в основном идет по двум направлениям: одни пакеты добавляют возможности управления чрезвычайно большими наборами данных, ускоряя операции с данными или допуская, чтобы размер данных превышал объем доступной системной памяти; другие позволяют R работать быстрее, в том числе распределяя работу между несколькими компьютерами или процессорами, используя специализированное компьютерное оборудование или предоставляя алгоритмы машинного обучения, оптимизированные для задач с большими объемами данных.

Управление сверхбольшими наборами данных

Сверхбольшие наборы данных могут привести к зависанию R из-за того, что системе не хватает памяти для хранения данных. Даже если весь набор

данных помещается в доступной памяти, необходима дополнительная память для обработки этих данных. Кроме того, обработка очень больших наборов данных, возможно, потребует много времени просто из-за большого объема записей; даже быстрая операция может привести к сбою, если ее необходимо выполнить несколько миллионов раз.

Пару лет назад предлагалось выполнить подготовку данных вне R на другом языке программирования, а затем использовать R для анализа небольшого подмножества данных. Однако в этом больше нет необходимости, поскольку в R появилось несколько пакетов для решения проблем, связанных с большими данными.

Быстрое построение фреймов данных с помощью `data.table`

Пакет `data.table`, разработанный Мэттом Доулом (Matt Dowle), Томом Шортом (Tom Short), Стивом Лианоглу (Steve Lianoglou) и Аруном Сринивасаном (Arjun Srinivasan), предлагает расширенную версию фрейма данных, называемую *таблицей данных*. Объекты `data.table` обычно работают намного быстрее, чем фреймы данных, для операций разбиения на подмножества, объединения и группировки. А для сверхбольших наборов данных, насчитывающих миллионы строк, такие объекты могут быть значительно быстрее, чем даже объекты `dplyr`. Однако, поскольку объект `data.table` по существу является улучшенным фреймом данных, полученные в результате объекты могут по-прежнему использоваться любой R-функцией, которая принимает фреймы данных.



Проект `data.table` размещен в GitHub по адресу <https://github.com/Rdatatable/data.table/wiki>.

После установки пакета `data.table` функция `fread()` будет считывать табличные файлы, такие как CSV, в объекты таблиц данных. Например, для того, чтобы загрузить использованные ранее данные о кредитах, введите следующее:

```
> library(data.table) > credit <-  
fread("credit.csv")
```

Затем можно запросить таблицу данных о кредитовании, используя синтаксис, подобный принятой в R форме `[row, col]`, — результат будет оптимизирован для обеспечения высокой скорости и некоторых дополнительных полезных свойств. В частности, структура таблицы данных позволяет указывать вместо строки сокращенную команду выбора строк, а вместо столбца — функцию, которая что-то делает с выбранными строками. Например, следующая команда вычисляет среднюю сумму кредита, которую запрашивают клиенты с хорошей кредитной историей:

```
> credit[credit_history == "good",  
mean(amount)][1] 3040.958
```

Создавая большие запросы с подобным простым синтаксисом, с таблицами данных можно выполнять очень сложные операции. А поскольку эта структура данных оптимизирована по скорости, ее можно использовать для больших наборов данных.

Единственное ограничение структур `data.table` состоит в том, что они, подобно фреймам данных, ограничены доступной системной памятью. В следующих двух разделах будут рассмотрены пакеты, которые устраняют этот недостаток за счет нарушения совместимости с многими R-функциями.



Пакеты `dplyr` и `data.table` имеют уникальные преимущества. Их подробное сравнение вы найдете в дискуссии Reddit по адресу https://www.reddit.com/r/rstats/comments/acjr9d/dplyr_performance/ и в аналогичном обсуждении на Stack Overflow: <https://stackoverflow.com/questions/21435339/data-table-vs-dplyr-can-one-do-something-well-the-other-cant-or-does-poorly>. Можно также извлечь лучшее из обоих пакетов, поскольку структуры `data.table` можно загружать в `dplyr` с помощью функции `tbl_dt()`.

Размещение фреймов данных на диске с помощью пакета `ff`

Пакет `ff`, созданный Даниэлем Адлером (Daniel Adler), Кристианом Глезером (Christian Glaeser), Олегом Ненадичем (Oleg Nenadic), Йенсом Эльшлегелем (Jens Oehlschlägel) и Вальтером Цуккини (Walter Zucchini), предоставляет альтернативу фрейму данных — объект `ffdf`, который позволяет создавать наборы данных из нескольких миллиардов строк, даже если это намного превышает доступную системную память.

Структура `ffdf` состоит из двух компонентов: физического, который хранит данные на диске в высокоэффективной форме, и виртуального, который внешне действует как типичный фрейм данных R, но при этом ссылается на данные, хранящиеся на физическом компоненте. Объект `ffdf` можно представить как отображение, которое указывает на местоположение данных на диске.



Проект `ff` размещен в Интернете по адресу <http://ff.r-forge.r-project.org/>. Недостатком структур данных `ffdf` является то, что большинство R-функций не может использовать их в чистом виде. Необходимо обрабатывать эти данные небольшими порциями и затем объединять

результаты. Преимущество разделения данных состоит в том, что задачу можно разделить между несколькими процессорами и выполнять одновременно, используя методы параллельных вычислений, представленные далее в этой главе.

После установки пакета `ff` используйте функцию `read.csv.ffdf()`, чтобы прочитать большой CSV-файл:

```
> library(ff)> credit <- read.csv.ffdf(file =  
"credit.csv", header = TRUE)
```

К сожалению, мы не можем работать с объектом `ffdf` напрямую, поскольку попытка обработать его как обычный фрейм данных приводит к появлению сообщения об ошибке:

```
> mean(credit$amount)[1] NA  
Warning message: In mean.default(credit$amount) :  
argument is not numeric or logical: returning NA
```

Пакет `ffbase`, созданный Эдвином де Йонгом (Edwin de Jonge), Яном Вейфелсом (Jan Wijffels) и Яном ван дер Лааном (Jan van der Laan), несколько упрощает эту проблему, добавляя возможности базового анализа для объектов `ff`. Это позволяет использовать объекты `ff` непосредственно для исследования данных. Например, после установки пакета `ffbase` функция `mean` работает, как предполагалось:

```
> library(ffbase)> mean(credit$amount)[1]  
3271.258
```

Этот пакет предоставляет и другие базовые функции, такие как математические операторы, функции запросов, сводную статистику и оболочки для работы с оптимизированными алгоритмами машинного обучения, такими как `biglm` (описан далее в этой главе). Эти меры не полностью устраняют проблемы работы со сверхбольшими наборами данных, но делают процесс более прозрачным.



Подробнее о расширенных функциональных возможностях `ffbase` читайте на сайте проекта по адресу <http://github.com/edwindj/ffbase>.

Использование больших матриц с помощью пакета `bigmemory`
Пакет `bigmemory`, созданный Майклом Дж. Кейном (Michael J. Kane), Джоном У. Эмерсоном (John W. Emerson) и Питером Хэверти (Peter Navearty), позволяет использовать очень большие матрицы, размер которых превышает объем доступной системной памяти. Матрицы могут храниться на диске или в общей памяти, что позволяет использовать их другими процессами на том же компьютере или по сети. Это облегчает параллельные вычисления, которые будут описаны далее в этой главе.



Подробную документацию по пакету `bigmemory` вы найдете по адресу <http://www.bigmemory.org/>.

Поскольку матрицы `bigmemory` намеренно созданы непохожими на фреймы данных, их нельзя использовать напрямую в большинстве методов ML, описанных в этой книге. Кроме того, в этих матрицах могут храниться только числовые данные. Тем не менее, поскольку матрицы `bigmemory` похожи на обычные R-матрицы, из них легко создавать меньшие выборки или порции и затем преобразовывать в стандартные структуры данных R.

Авторы пакета `bigmemory` также являются разработчиками пакетов `bigalgebra`, `biganalytics` и `bigtabulate`, которые позволяют выполнять простой анализ матриц. Особо следует отметить функцию `bigkmeans()` из пакета `biganalytics`, которая выполняет кластеризацию методом k-средних (описан в главе 9). Вследствие узкоспециализированного характера этих пакетов примеры их использования выходят за рамки данной главы.

Ускорение обучения благодаря параллельным вычислениям

На заре компьютерной техники процессоры выполняли инструкции последовательно. Это означало, что они ограничивались выполнением только одной операции за раз. Следующая инструкция не могла быть запущена, пока не завершится выполнение предыдущей (рис. 12.5). Уже тогда было понятно, что многие задачи можно решать более эффективно, выполняя несколько этапов одновременно, однако необходимая для этого технология еще попросту не существовала.

Последовательные вычисления

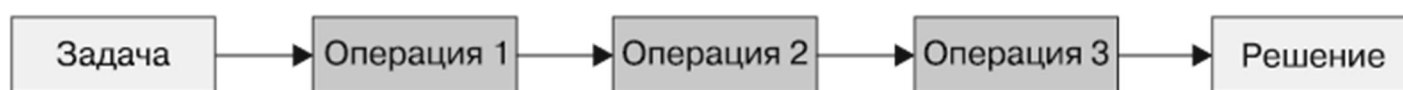


Рис. 12.5. При последовательных вычислениях невозможно начать следующую операцию, не закончив выполнение предыдущей

Эта проблема была решена путем разработки методов *параллельных вычислений*, в которых для решения больших задач используются комплексы из двух и более процессоров или компьютеров. Многие современные компьютеры изначально рассчитаны на параллельные вычисления. Даже если у них только один процессор, у него обычно несколько ядер, способных работать параллельно. Это позволяет выполнять задачи независимо друг от друга (рис. 12.6).

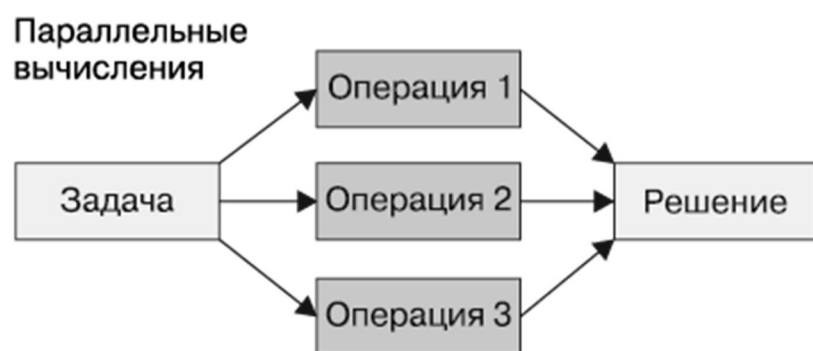


Рис. 12.6. При параллельных вычислениях операции выполняются одновременно. В конце их результаты объединяются

Сети, составленные из нескольких компьютеров — кластеры, — также можно использовать для параллельных вычислений. Большой кластер может включать в себя различные аппаратные средства, которые могут находиться на значительном расстоянии друг от друга. В этом случае кластер называется грид-системой (grid). Численность компьютеров в кластере или грид-системе может достигать сотен или тысяч машин. Даже если это обычное оборудование, вместе оно может образовывать очень мощные системы.

Однако подвох в том, что не любую задачу можно распараллелить. Одни задачи распараллеливаются лучше, другие — хуже. Может показаться, что добавление в систему 100 процессоров приведет к ускорению работы в 100 раз (то есть общее время выполнения составит $1/100$), но обычно это не так. Причина в том, что для управления этими компьютерами требуются дополнительные усилия. Работа должна быть разделена на равные непересекающиеся задачи, и в конце все результаты, полученные с отдельных машин, должны быть объединены в один окончательный ответ (см. рис. 12.6).

Так называемые *легко распараллеливаемые* задачи являются идеальным случаем. Эти задачи легко сводятся к непересекающимся рабочим блокам, а их результаты легко объединяются. Примером легко распараллеливаемой задачи машинного обучения является десятиблочная кросс-валидация; после разделения данных на десять выборок каждый из десяти блоков становится независимым — другими словами, не влияет на остальные блоки. Как вы скоро увидите, решение такой задачи можно значительно ускорить с помощью параллельных вычислений.

Измерение времени выполнения

Все усилия по ускорению R были бы напрасными, если бы невозможно было систематически измерять количество сэкономленного времени. Одним из вариантов является секундомер, однако более простое решение — обернуть интересующий нас код в функцию `system.time()`.

Например, на моем ноутбуке функция `system.time()` сообщает, что генерация миллиона случайных чисел занимает около 0,08 секунды:

```
>
```

```
system.time(rnorm(1000000)) user system elapsed  
0.079 0.000 0.067
```

Эта же функция может использоваться для оценки улучшения производительности, достигнутой с помощью только что описанных методов или любой другой R-функции.



Вообще-то, когда вышло первое издание этой книги, генерация миллиона случайных чисел занимала 0,130 секунды; для второго издания та же операция заняла около 0,093 секунды. Теперь она занимает около 0,067 секунды. Хотя каждый раз я использовал немного более мощный компьютер, это ускорение примерно на 50 % в течение примерно шести лет показывает, насколько быстро совершенствуются компьютерное оборудование и программное обеспечение.

Параллельные вычисления с помощью пакетов `multicore` и `snow` Пакет `parallel`, включенный в R 2.14.0 и более поздние версии, упростил развертывание параллельных алгоритмов, предоставляя стандартную структуру для настройки рабочих процессов, позволяющих выполнять задачи одновременно. Это достигается путем включения компонентов из пакетов `multicore` и `snow`, каждый из которых использует свой подход к обеспечению многозадачности.

Если ваш компьютер достаточно новый, то вы, вероятно, сможете реализовать на нем параллельную обработку данных. Чтобы определить количество ядер на компьютере, используйте функцию `detectCores()`. Обратите внимание, что ваш результат будет отличаться в зависимости от характеристик оборудования:

```
> library(parallel)> detectCores()[1] 8
```

Пакет `multicore`, разработанный Саймоном Урбанеком (Simon Urbanek), позволяет выполнять параллельную обработку данных на одном компьютере с несколькими процессорами или процессорными ядрами. Этот пакет использует возможности многозадачности, заложенные в операционной системе компьютера, для создания дополнительных R-сессий, которые совместно используют одну и ту же память. Возможно, этот пакет — самый простой способ начать параллельную обработку данных в R. К сожалению, поскольку Windows не поддерживает параллельные вычисления, в этой операционной системе данное решение не будет работать.

Простой способ начать работу с многоядерными функциями — использовать функцию `mclapply()`, которая является многоядерной

версией `lapply()`. Например, следующие блоки кода иллюстрируют, как распределить задачу генерации миллиона случайных чисел между 1, 2, 4 и 8 ядрами. После того как все ядра завершат свою работу, для объединения параллельных результатов (списка) в один вектор используется функция `unlist()`:

```
> system.time(l1 <- unlist(mclapply(1:10,
function(x) {+ rnorm(1000000)}), mc.cores =
1))) user system elapsed 0.627 0.015 0.
647> system.time(l2 <- unlist(mclapply(1:10,
function(x) {+ rnorm(1000000)}), mc.cores =
2))) user system elapsed 0.751 0.211 0.
568> system.time(l4 <- unlist(mclapply(1:10,
function(x) {+ rnorm(1000000)}), mc.cores =
4))) user system elapsed 0.786 0.270 0.
405> system.time(l8 <- unlist(mclapply(1:10,
function(x) {+ rnorm(1000000)}), mc.cores =
8))) user system elapsed 1.033 0.315 0.
321
```

Обратите внимание, как с увеличением количества ядер уменьшается время выполнения, и каждый раз на меньшую величину. Это простой пример, однако его можно легко адаптировать ко многим другим задачам.

Пакет `snow` (от *Simple Networking of Workstations* — «простая сеть из рабочих станций») разработан Люком Тирни (Luke Tierney), Э. Дж. Россини (A. J. Rossini), На Ли (Na Li) и Х. Севчиковой (H. Sevcikova). Этот пакет позволяет выполнять параллельные вычисления на многоядерных и многопроцессорных компьютерах, а также в сети, состоящей из нескольких компьютеров. Он немного сложнее в использовании, чем `multicore`, но гораздо мощнее и гибче.

Функциональность `snow` включена в пакет `parallel`, поэтому для настройки кластера на одном компьютере используйте функцию `makeCluster()`, указав количество используемых ядер:

```
> cl1 <- makeCluster(4)
```

Поскольку `snow` включает данные в сетевой трафик, вы, в зависимости от используемой операционной системы, можете получить сообщение, требующее подтверждения доступа через брандмауэр.

Чтобы убедиться в работоспособности кластера, можно попросить каждый узел сообщить его имя хоста.

Функция `clusterCall()` выполняет функцию на каждом компьютере кластера. В данном случае мы определим функцию, которая просто вызывает функцию `Sys.info()` и возвращает параметр `nodename`:


```

> clusterCall(cl1, function() {
Sys.info()["nodename"] }
)[[1]]          nodename "Bretts-Macbook-
Pro.local" [[2]]          nodename "Bretts-
Macbook-
Pro.local" [[3]]          nodename "Bretts-
Macbook-
Pro.local" [[4]]          nodename "Bretts-
Macbook-Pro.local"

```

Поскольку все четыре узла работают на одной машине, неудивительно, что они сообщают одно и то же имя хоста. Чтобы на четырех узлах выполнялись разные команды, предоставьте каждому из них уникальный параметр с помощью функции `clusterApply()`. Здесь мы передадим узлам разные буквы. Затем узлы будут параллельно выполнять простую функцию, каждый для своей буквы:

```

> clusterApply(cl1, c('A', 'B', 'C',
'D'), function(x) {
paste("Cluster", x, "ready!") })[[1]][1] "Cluster
A ready!" [[2]][1] "Cluster B ready!" [[3]][1]
"Cluster C ready!" [[4]][1] "Cluster D ready!"

```

Заканчивая работу с кластером, важно завершить все процессы, которые он начал, чтобы освободить ресурсы, используемые узлами кластера:

```

> stopCluster(cl1)

```

С помощью этих простых команд можно ускорить выполнение многих задач машинного обучения. Для наиболее серьезных проблем, связанных с большими данными, существуют более сложные конфигурации `snow`. Например, можно попробовать построить *кластер класса* Beowulf — сеть, состоящую из большого количества персональных компьютеров. Для академических и отраслевых исследований с выделенными вычислительными кластерами пакет `snow` может использовать пакет `Rmpi` для доступа к высокопроизводительным серверам интерфейса передачи сообщений (Message Passing Interface, MPI). Работа с такими кластерами требует знания сетевых конфигураций и вычислительного оборудования. В данной книге информация об этом не приводится.



Более подробное введение в `snow`, включая некоторую информацию о том, как настроить параллельные вычисления на нескольких компьютерах по сети, вы найдете в лекции Люка Тирни (Luke Tierney) по адресу <http://homepage.stat.uiowa.edu/~luke/classes/295-hpc/notes/snow.pdf>.

Использование преимуществ параллельных вычислений с помощью пакетов `foreach` и `doParallel`

Пожалуй, начать использовать параллельные вычисления будет проще всего с пакета `foreach`, созданного Ричем Калауэем (Rich Calaway) и Стивом Уэстоном (Steve Weston), особенно если вы работаете с R в операционной системе Windows, поскольку многие другие пакеты зависят от платформы.

Ядром пакета `foreach` является конструкция цикла `foreach`. Если вы работали с другими языками программирования, то вам она должна быть знакома. По сути, эта конструкция позволяет перебирать в цикле элементы множества, не вычисляя их количество; другими словами, она означает «**сделать что-нибудь для каждого** элемента из множества данных».

Если вы думаете, что в R и так уже есть набор функций для циклического перебора множества элементов (например, `apply()`, `lapply()`, `sapply()` и т.д.), то вы правы. Однако у цикла `foreach` есть дополнительное преимущество: итерации цикла можно распараллелить с использованием очень простого синтаксиса. Посмотрим, как это работает.

Вспомним команду, которую мы использовали для генерации миллиона случайных чисел. Усложним задачу: увеличим количество чисел до 100 000 000, в результате чего процесс займет более шести секунд:

```
> system.time(l1 <-  
rnorm(100000000)) user system elapsed 5.873  
0.204 6.087
```

После установки пакета `foreach` эту задачу можно описать как цикл, который генерирует четыре набора по 25 000 000 случайных чисел. Параметр `.combine` является необязательным; он сообщает `foreach`, какую функцию следует использовать для объединения набора результатов, полученных на каждой итерации цикла, в окончательный результат. В данном случае, поскольку каждая итерация генерирует набор случайных чисел, мы просто воспользуемся функцией конкатенации `c()` и создадим общий объединенный вектор:

```
> library(foreach)> system.time(l4 <- foreach(i  
= 1:4, .combine = 'c') %do%  
rnorm(25000000) user system elapsed 6.177  
0.391 6.578
```

Вы заметили, что эта функция не привела к улучшению скорости? Отлично! На самом деле процесс был даже медленнее. Причина в том, что по умолчанию пакет `foreach` запускает итерации цикла последовательно,

а функция добавляет в процесс еще некоторое количество вычислительных затрат. Для того чтобы добавить `foreach` параллелизм, нужно воспользоваться смежным пакетом `doParallel`, который, в свою очередь, использует пакет `parallel`, описанный ранее в этой главе и входящий в состав R. После установки пакета `doParallel` зарегистрируйте количество ядер и замените команду `%do%` на `%dopar%`:

```
> library(doParallel)> registerDoParallel(cores
= 4)> system.time(l4p <- foreach(i = 1:4,
.combine = 'c')
rnorm(25000000)) user system elapsed 7.841
2.288 3.894
```

Как видно по результатам, это привело к ожидаемому увеличению производительности, время выполнения сократилось примерно на 40 %.

Чтобы закрыть кластер `doParallel`, просто введите команду:

```
> stopImplicitCluster()
```

По окончании сеанса R кластер закрывается автоматически, но лучше сделать это напрямую.

Обучение и оценка эффективности параллельных моделей с помощью пакета `caret`

Пакет `caret` Макса Куна (подробно рассмотренный в главах 10 и 11) использует параллельные пакеты, если они зарегистрированы в R с использованием описанного ранее пакета `foreach`.

Рассмотрим простой пример, в котором обучим модель случайного леса на наборе данных о кредитах. Без распараллеливания обучение модели занимает около 79 секунд:

```
> library(caret)> credit <-
read.csv("credit.csv")> system.time(train(default
~ ., data = credit, method =
"rf")) user system elapsed77.345 1.778 79
.205
```

Если же использовать пакет `doParallel` и зарегистрировать четыре ядра для параллельной работы, то на построение модели уходит менее 20 секунд — меньше четверти исходного времени — и нам не потребовалось менять ни одну строку кода `caret`:

```
> library(doParallel)> registerDoParallel(cores
= 8)> system.time(train(default ~ ., data =
credit, method =
```

```
"rf" ) ) user system elapsed122.579 3.292 1
9.034
```

Многие задачи, связанные с обучением и оценкой эффективности моделей, такие как создание случайных выборок и многократное тестирование прогнозов для десятиблочной кросс-валидации, легко распараллеливаются и давно созрели для повышения эффективности. С учетом этого целесообразно всегда регистрировать несколько ядер перед началом проекта `caret`.



На веб-сайте проекта `caret` вы найдете более подробные инструкции по настройке и пример повышения эффективности при использовании параллельной обработки в `caret`: <http://topepo.github.io/caret/parallel.html>.

Облачные параллельные вычисления с помощью MapReduce и Hadoop

Программная модель *MapReduce* была разработана в Google как способ обработки данных на большом кластере компьютеров, объединенных в сеть. Модель MapReduce определяет параллельное программирование как процесс, состоящий из следующих двух этапов:

- этапа *map*, на котором исходная задача делится на более мелкие задачи, затем распределяемые между компьютерами кластера;
- этапа *reduce*, на котором результаты этих более мелких кусков работы собираются и синтезируются в окончательное решение исходной задачи.

Популярной альтернативой закрытой среде MapReduce является ApacheHadoop с открытым исходным кодом. Программное обеспечение Hadoop включает в себя концепцию MapReduce и распределенную файловую систему, способную хранить большие объемы данных на кластере компьютеров.



Издательство Packt Publishing опубликовало большое количество книг по Hadoop. Последние издания вы найдете на странице <https://www.packtpub.com/all/?search=hadoop>.

В настоящее время разрабатываются несколько проектов R, которые предоставляют R-интерфейс для Hadoop. Один из них — проект RNadoop от Revolution Analytics. Этот проект предоставляет пакет `rnr2`, предназначенный для того, чтобы разработчики в среде R могли легко писать программы для MapReduce. Еще один аналогичный пакет, `plyrnr`, обеспечивает функциональность, аналогичную пакету `dplyr`, для обработки больших наборов данных. В состав RNadoop

также входят пакеты, предоставляющие R-функции для доступа к распределенным хранилищам данных Hadoop.



Подробнее о проекте RHadoop читайте

здесь: <https://github.com/RevolutionAnalytics/RHadoop/wiki>.

Несмотря на то что Hadoop является достаточно зрелым фреймворком, для его использования и выполнения даже простейших задач машинного обучения необходимы специальные навыки программирования. Возможно, именно этим объясняется очевидное отсутствие популярности этого пакета у пользователей R. Кроме того, хотя Hadoop отлично работает со сверхбольшими объемами данных, этот проект не всегда является самым быстрым вариантом, поскольку он хранит все данные на диске, вместо того чтобы использовать доступную память. В следующем разделе мы рассмотрим расширение Hadoop, которое решает эти проблемы скорости и удобства использования.

Параллельные облачные вычисления с помощью Apache Spark

Проект Apache Spark — это кластерная вычислительная среда для обработки больших данных, имеющая множество преимуществ по сравнению с Apache Hadoop. Поскольку Apache Spark использует доступную память кластера, эта среда обрабатывает данные примерно в 100 раз быстрее, чем Hadoop. Кроме того, она предоставляет высокоуровневые библиотеки для многих распространенных задач обработки, анализа и моделирования данных. Сюда относятся язык запросов данных SparkSQL, библиотека машинного обучения MLlib, GraphX для анализа графов и сетей и библиотека Spark Streaming для обработки потоков данных в реальном времени.



Издательство Packt Publishing выпустило множество книг по Spark.

Последние издания вы найдете на

странице <https://www.packtpub.com/all/?search=spark>.

Apache Spark часто выполняется удаленно на кластере виртуальных машин, размещенном в облаке, но преимущества этой модели можно увидеть и на локальном оборудовании. В любом случае

пакет `sparklyr` подключается к кластеру и предоставляет

интерфейс `dplyr` для анализа данных с использованием Spark. Более подробные инструкции по использованию Spark в среде R вы найдете по адресу <https://spark.rstudio.com>, но основы, необходимые для начала работы, весьма просты.

Чтобы проиллюстрировать основные принципы, мы построим модель случайного леса на основе набора данных о кредитах для прогнозирования невозврата кредита. Начнем с установки пакета `sparklyr`. Затем мы создадим экземпляр кластера Spark на локальном компьютере, используя следующий код:

```
> install.packages("sparklyr")>
library(sparklyr)> spark_install(version =
"2.1.0")> spark_cluster <- spark_connect(master =
"local")
```

Затем мы загрузим набор данных о кредитах из файла `credit.csv` на наш локальный компьютер в экземпляр Spark, а затем используем функцию разделения фрейма данных Spark `sdf_partition()`, чтобы случайным образом разделить данные на тренировочный и тестовый наборы в пропорции 75 и 25 % соответственно. Параметр `seed` — это случайное начальное число, обеспечивающее идентичность результатов при каждом запуске кода:

```
> credit_spark <- spark_read_csv(spark_cluster,
"credit.csv")> splits <-
sdf_partition(credit_spark, train =
0.75, test = 0.25, seed
= 123)
```

После этого мы передадим тренировочные данные в функцию модели случайного леса, сделаем прогнозы и используем оценку классификации для вычисления AUC на тестовом наборе данных:

```
> credit_rf <- splits$train
%>% ml_random_forest(default ~ .)> pred <-
ml_predict(credit_rf, splits$test)>
ml_binary_classification_evaluator(pred, metri
c_name = "areaUnderROC")[1] 0.7848068
```

С помощью всего лишь нескольких строк R-кода мы создали модель случайного леса с использованием Spark. Эту модель можно расширять для моделирования миллионов записей. Если потребуется еще больше вычислительной мощности, код можно будет выполнить в облаке с использованием кластера интенсивных параллельных вычислений Spark, просто передав функции `spark_connect()` соответствующее имя хоста. Код также легко адаптируется к другим методам моделирования с помощью одной из функций обучения с учителем, полный список которых находится по адресу <https://spark.rstudio.com/mlib/>.



Вероятно, самый простой способ начать работу со Spark — воспользоваться облачной платформой Databricks, разработанной создателями Spark. Веб-интерфейс Databricks позволяет легко управлять кластерами и масштабировать их. Бесплатная версия community edition, которая находится на сайте <https://databricks.com>, предоставляет небольшой кластер, чтобы можно было попрактиковаться на нескольких обучающих материалах или даже поэкспериментировать с собственными данными.

Развертывание оптимизированных алгоритмов обучения

Некоторые алгоритмы машинного обучения, описанные в этой книге, способны обрабатывать очень большие наборы данных с относительно незначительными изменениями. Например, было бы довольно легко реализовать наивный байесовский алгоритм или алгоритм Apriori, используя одну из описанных в предыдущих разделах структур, предназначенных для сверхбольших наборов данных. Некоторые типы алгоритмов обучения, такие как ансамбли, хорошо поддаются распараллеливанию, поскольку работа каждой модели может быть распределена между процессорами или компьютерами кластера. Другие модели, напротив, требуют больших изменений данных или алгоритма, или же их необходимо полностью переосмыслить, прежде чем удастся применить к сверхбольшим наборам данных.

В следующих разделах будут описаны пакеты, в которых представлены оптимизированные версии алгоритмов обучения, рассмотренных ранее.

Построение больших регрессионных моделей с помощью пакета `biglm`

Пакет `biglm`, разработанный Томасом Ламли (Thomas Lumley), предоставляет функции для обучения регрессионных моделей на наборах данных, которые могут оказаться слишком большими и не поместиться в памяти. Модель обновляется итеративно, используя небольшие порции данных. Несмотря на измененный подход, результаты оказываются практически такими же, как и те, что получены с использованием обычной функции `lm()` для всего набора данных.

Для удобства работы со сверхбольшими наборами данных функция `biglm()` позволяет заменить фрейм данных базой данных SQL. Модель также может быть обучена на порциях данных, полученных из объектов данных, созданных с помощью описанного ранее пакета `ff`.

Ускоренное выращивание случайного леса с помощью пакета `ranger`

Пакет `ranger` Марвина Н. Райта (Marvin N. Wright), Стефана Вагера (Stefan Wager) и Филиппа Пробста (Philipp Probst) выполняет ускоренный вариант алгоритма случайного леса, специально предназначенный для

наборов данных с большим количеством признаков или примеров.

Функция используется почти так же, как описанный ранее случайный лес:

```
> library(ranger)> credit <-  
read.csv("credit.csv")> m <- ranger(default ~ .,  
data = credit, num.trees =  
500, mtry = 4)> p <- predict(m,  
credit)
```

Обратите внимание, что, в отличие от большинства представленных ранее результатов функции `predict()`, прогнозы `ranger` хранятся в виде подобъекта в объекте `predictions`:

```
> head(p$predictions)[1]  
no yes no no yes noLevels: no yes
```

Использование функции `ranger()` — самый простой способ построить расширенный и улучшенный случайный лес, не прибегая к кластерным вычислениям или альтернативным структурам данных.

Выращивание сверхбольших случайных лесов с помощью пакета `bigrf`

В пакете `bigrf` Алоизиуса Лима (Aloysius Lim) реализовано обучение случайного леса для классификации и регрессии на наборах данных, размер которых слишком велик, чтобы поместиться в память. Для решения этой проблемы в пакете `bigrf` используются объекты `bigmemory`, описанные ранее в этой главе. Для более быстрого построения леса можно использовать этот пакет совместно с описанными ранее пакетами `foreach` и `doParallel`, чтобы реализовать параллельный рост деревьев.



Подробнее о пакете `bigrf`, включая примеры и инструкции по установке в Windows, читайте на вики-странице пакета по адресу <https://github.com/alloysius-lim/bigrf>.

Процессор ускоренного машинного обучения H2O

Проект H2O — это инфраструктура больших данных, в которой предоставлены ускоренные реализации алгоритмов машинного обучения в оперативной памяти. Эти реализации также могут работать в среде кластерных вычислений. Проект включает в себя функции для многих методов, описанных в этой книге, включая наивный байесовский алгоритм, регрессию, глубокие нейронные сети, кластеризацию методом k-средних, ансамблевые методы и случайные леса, а также многие другие алгоритмы.

H2O использует эвристики для нахождения приближенных решений в задачах машинного обучения, повторно итерируя над небольшими порциями данных. Это позволяет точно определять, какое количество данных из сверхбольшого набора следует использовать для обучения алгоритма. Для одних задач подходит более быстрое решение, для других может потребоваться полный набор, что займет дополнительное время для обучения.

H2O обычно работает со сверхбольшими наборами данных значительно быстрее и лучше, чем функции машинного обучения Apache Spark (MLlib), которые зачастую намного быстрее, чем базовые функции R. Однако, поскольку Apache Spark является широко используемой средой для кластерных вычислений и подготовки больших данных, H2O можно запустить на Apache Spark с помощью программного обеспечения *Sparkling Water*. Благодаря Sparkling Water ученые-исследователи получают лучшее из обоих проектов: преимущества Spark для подготовки данных и преимущества H2O для машинного обучения.

Пакет `h2o` предоставляет функциональные возможности для доступа к экземпляру H2O из среды R. Полное руководство по H2O не приводится в этой книге, всю документацию по проекту вы найдете по адресу <http://docs.h2o.ai>. Отмечу лишь, что основы H2O просты. Для начала установите и загрузите пакет `h2o`. Затем инициализируйте локальный экземпляр H2O, используя следующий код:

```
> library(h2o) > h2o_instance <- h2o.init()
```

В результате на вашем компьютере запустится сервер H2O, который можно просмотреть через веб-приложение H2O Flow по адресу <http://localhost:54321>. H2O Flow позволяет администрировать и отправлять команды на сервер H2O, и даже создавать модели и оценивать их эффективность с помощью простого браузерного интерфейса (рис. 12.7).

Мы могли бы выполнить анализ в этом интерфейсе, но давайте вернемся к R и используем H2O на данных о кредитах, которые мы уже изучали ранее. Для этого сначала нужно загрузить набор данных `credit.csv` в этот экземпляр с помощью следующей команды:

```
> credit.hex <- h2o.uploadFile("credit.csv")
```

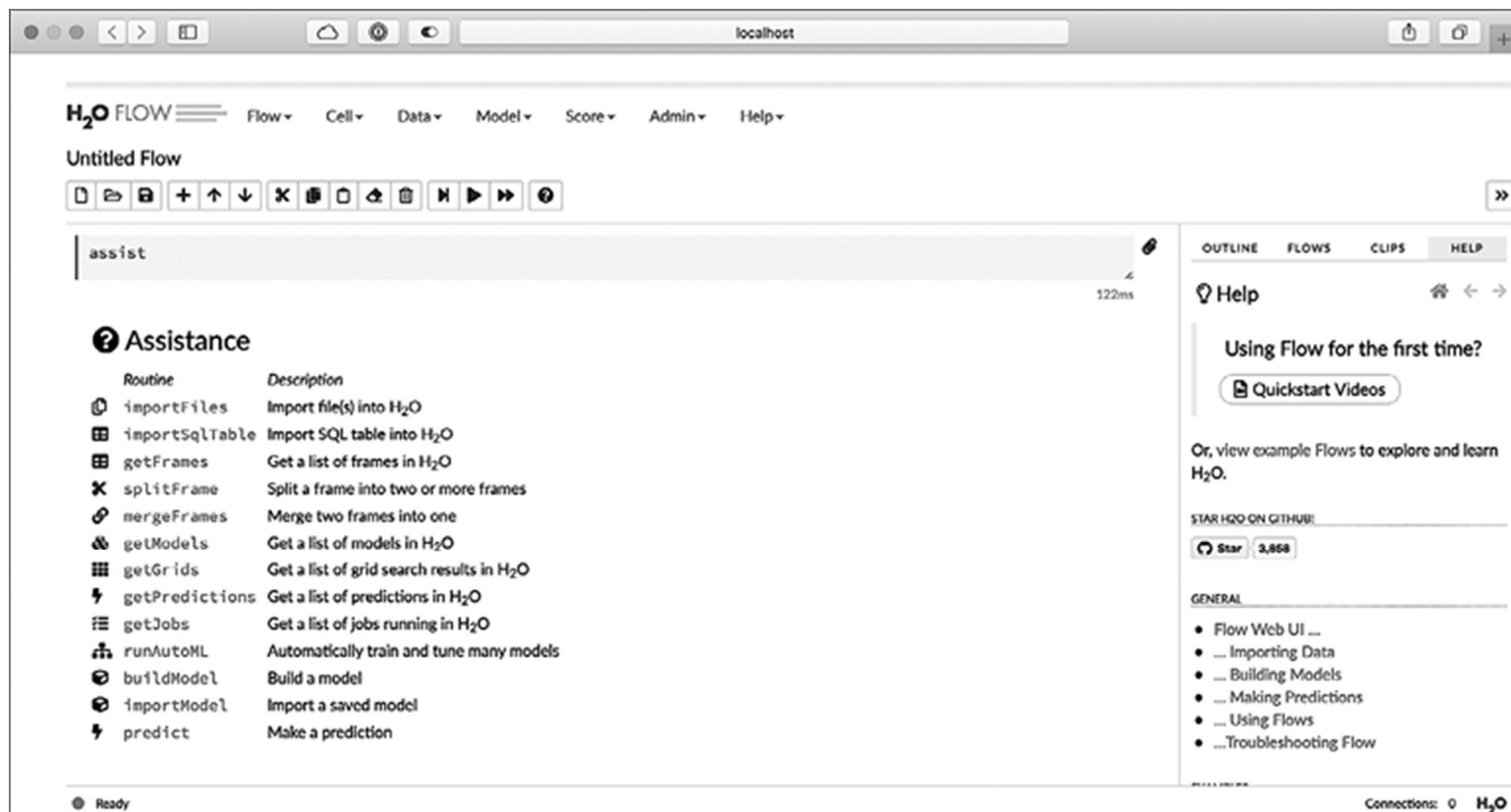


Рис. 12.7. H2O Flow — веб-приложение для взаимодействия с экземпляром H2O

Обратите внимание, что фрейм данных H2O обозначается с помощью расширения `.hex`.

Далее мы применим к этому набору данных реализацию случайного леса H2O с помощью следующей команды:

```
> h2o.randomForest(y =
"default", training_frame =
credit.hex, ntrees =
500, seed = 123)
```

Результаты работы этой команды включают в себя информацию об оценках эффективности исходной модели:

```
** Reported on training data. *** Metrics
reported on Out-Of-Bag training samples
**MSE:      0.1636964RMSE:      0.4045941LogLoss:
0.4956524Mean Per-Class
Error:      0.2835714AUC:      0.7844881pr_auc:      0.
6192192Gini:      0.5689762
```

Использованный здесь набор данных о кредитах не очень велик, однако этот код H2O можно масштабировать для наборов данных практически любого размера. Кроме того, код практически не изменится, если понадобится запустить его в облаке — просто передайте функции `h2o.init()` адрес удаленного хоста.

Вычисления на GPU

Альтернативой параллельной обработке данных является использование графического процессора (GPU) компьютера для увеличения скорости математических вычислений. Графический процессор — это специализированный процессор, который оптимизирован для быстрого

вывода изображений на экран компьютера. Поскольку компьютеру часто приходится отображать сложную трехмерную графику (особенно для видеоигр), аппаратное обеспечение многих графических процессоров рассчитано на параллельную обработку и оптимизировано для чрезвычайно эффективных матричных и векторных вычислений. Дополнительным преимуществом графических процессоров является то, что их можно использовать для эффективного решения определенных типов математических задач. Если обычный процессор ноутбука или настольного компьютера может иметь до 16 ядер, то в среднем графическом процессоре количество ядер может достигать нескольких тысяч или даже десятков тысяч (рис. 12.8).

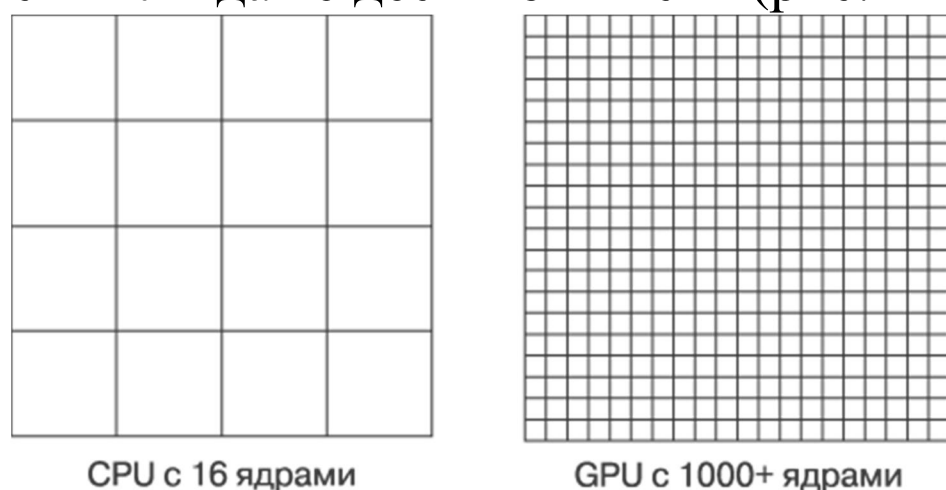


Рис. 12.8. В графическом процессоре (GPU) во много раз больше ядер, чем в обычном центральном процессоре компьютера (CPU)

Недостатком вычислений на GPU является то, что для этого требуется специальное оборудование, которое не входит в комплектацию многих компьютеров. В большинстве случаев требуется графический процессор NVIDIA — в нем реализована закрытая среда разработки *Complete Unified Device Architecture (CUDA)*, которая позволяет писать программы для GPU на обычных языках программирования, таких как C++.



Подробнее о роли NVIDIA в вычислениях на GPU читайте на странице <https://www.nvidia.com/en-us/about-nvidia/ai-computing/>.

В пакете `gputools` Джоша Бакнера (Josh Buckner), Марка Селигмана (Mark Seligman) и Джастина Уилсона (Justin Wilson) реализовано несколько R-функций, в том числе матричные операции, кластеризация и регрессионное моделирование, с использованием инструментария NVIDIA CUDA. Для работы этого пакета требуется графический процессор CUDA версии 1.3 или выше и установленный на компьютере инструментарий NVIDIA CUDA.

Гибкие числовые вычисления и машинное обучение с помощью
TensorFlow

Одним из наиболее значительных недавних инноваций в программном обеспечении машинного обучения является *TensorFlow* (<https://www.tensorflow.org>) — математическая библиотека с открытым исходным кодом, разработанная Google для углубленного ML. TensorFlow предоставляет вычислительный интерфейс на основе ориентированных графов, которые пропускают (англ. *flow*) массивы данных, называемые *тензорами*, через ряд математических операций. Таким образом, даже очень сложные методы «черного ящика», такие как глубокие нейронные сети, можно представить в виде более простых абстракций. Кроме того, поскольку наборы вычислений хранятся в графах как наборы взаимосвязанных операций, TensorFlow позволяет распределять работу между ядрами CPU или GPU и использовать преимущества сред с интенсивными параллельными вычислениями.



Издательство Packt Publishing выпустило много книг по TensorFlow. Последние издания вы найдете на странице <https://www.packtpub.com/all/?search=tensorflow>.

Команда RStudio разработала несколько R-интерфейсов для TensorFlow. Основной API реализован в виде пакета `tensorflow`, а в пакете `tfestimators` представлены высокоуровневые функции машинного обучения. Обратите внимание, что метод направленных графов TensorFlow можно использовать для реализации множества моделей ML, включая те, что были описаны в этой книге. Однако для этого требуется глубокое понимание матричной математики, которая определяет каждую модель, поэтому TensorFlow в данной книге не описывается. Для получения дополнительной информации об этих пакетах и возможностях RStudio по взаимодействию с TensorFlow посетите веб-сайт <https://tensorflow.rstudio.com>.



Благодаря уникальному подходу TensorFlow к оценке эффективности моделей машинного обучения вы увидите различия в терминологии, используемой приверженцами TensorFlow даже для таких простых методов, как линейная регрессия. Вам встретятся такие фразы, как «функция стоимости», «градиентный спуск» и «оптимизация». Эта терминология подтверждает, что машинное обучение с TensorFlow во многом аналогично построению нейронной сети, которая находит наилучшее приближение к желаемой модели.

Keras — интерфейс для глубокого обучения

Поскольку вычислительная среда TensorFlow идеально подходит для построения глубоких нейронных сетей, для обеспечения более простого высокоуровневого интерфейса с этим широко используемым функционалом была разработана библиотека Keras (<https://keras.io>). Keras написана на Python и может использовать TensorFlow и аналогичные фреймворки в качестве «движка» для внутренних вычислений. С помощью Keras можно выполнять глубокое обучение всего за несколько строк кода даже в таких сложных случаях, как классификация изображений.



Издательство Packt Publishing выпустило множество книг и видеоматериалов для изучения Keras. Последние издания вы найдете на странице <https://www.packtpub.com/all/?search=keras>.

Пакет `keras`, разработанный генеральным директором и основателем RStudio Джозефом Аллером (J. J. Allaire), представляет собой R-интерфейс к Keras. Хотя для использования пакета `keras` нужно совсем немного кода, разработка полезных моделей глубокого обучения требует глубоких знаний нейронных сетей, а также знания TensorFlow и API Keras. Кроме того, для создания чего угодно, кроме простейших нейронных сетей, важно использовать графический процессор — без интенсивного распараллеливания, которое обеспечивает графический процессор, код просто никогда не завершится. Поэтому обучение `keras` в данной книге не рассматривается. Обратитесь к документации RStudio по адресу <https://keras.rstudio.com> или к книге *Deep Learning with R* (2018), авторами которой являются Франсуа Шолле (Francois Chollet) и Дж. Аллер — создатели Keras и пакета `keras`. Это лучшие источники для начала изучения данного инструментария.



На момент публикации этой книги типичный графический процессор, используемый для глубокого обучения, стоил несколько сотен долларов за модель начального уровня и 1000–3000 долларов — за более производительную модель среднего уровня. Высококачественные устройства могут стоить много тысяч долларов. Чтобы не вкладывать так много денег сразу, многие почасово арендуют серверное время у облачных провайдеров, таких как Amazon AWS и Microsoft Azure. Стоимость минимального графического процессора у них составляет примерно 1 доллар в час — просто не забудьте отключиться, когда закончите работу, иначе это может обойтись очень дорого! Команда RStudio также предоставляет информацию о своем предпочтительном хосте по адресу https://tensorflow.rstudio.com/tools/cloud_desktop_gpu.html.

Резюме

Надеюсь, вы не пожалели, что уделили время изучению машинного обучения! Исследования относительно неизведанных границ параллельных и распределенных вычислений продолжаются и обещают большие возможности для применения знаний, полученных в процессе изучения больших данных. А росту сообщества специалистов в области анализа данных способствует свободный язык программирования R с открытым исходным кодом, который обеспечивает легкий доступ, — нужно лишь быть готовыми учиться.

Темы, рассмотренные в этой и предыдущих главах, закладывают основу для понимания более углубленных методов ML. Теперь вы просто обязаны продолжать учиться и добавлять в свой арсенал новые инструменты. И не забывайте о правиле отсутствия бесплатных завтраков: ни один алгоритм не решает всех проблем, у каждого из них есть свои сильные и слабые стороны. В машинном обучении всегда будет присутствовать человеческий фактор — знания о предметной области и возможность сопоставлять выбранный алгоритм с поставленной задачей.

В ближайшие годы будет интересно наблюдать, как изменяется этот человеческий фактор, поскольку грань между машинным обучением и обучением людей размыта. Такие сервисы, как Amazon Mechanical Turk, предоставляют интеллектуальные ресурсы краудсорсинга — кластер людей, готовых в любой момент выполнять простые задачи. Возможно, однажды, подобно тому, как мы использовали компьютеры для выполнения задач, которые нелегко даются людям, компьютеры станут использовать людей, чтобы делать обратное. Какая интересная тема для размышлений!

1 Этот и другие аспекты обеспечения кибербезопасности систем с машинным обучением раскрыты в книге: *Уорр К.* «Надежность нейронных сетей: укрепляем устойчивость ИИ к обману». — СПб.: Питер, 2020.

2 Впрочем, поскольку алгоритмы машинного обучения основаны на формировании статистической, а не логической модели, зачастую человек не способен понять (выделить слово), хоть и может увидеть процесс обработки отдельной порции данных (единицы анализа).

Научный редактор *Н. Искра*

Переводчик *Е. Сандицкая*

Литературный редактор *А. Дубейко*

Корректоры *Н. Искра, Е. Павлович, Е. Рафалюк-Бузовская*

Бретт Ланц

Машинное обучение на R: экспертные техники для прогностического анализа. — СПб.: Питер, 2020.

ISBN 978-5-4461-1512-9