

O'REILLY®



# Data Science

## Наука о данных с нуля



Джоэл Грас

Джоэл Грас

# Data Science

## Наука о данных с нуля

Санкт-Петербург  
«БХВ-Петербург»

2017

УДК 004.6  
ББК 32.81  
Г77

**Грас Дж.**

Г77 Data Science. Наука о данных с нуля: Пер. с англ. — СПб.: БХВ-Петербург, 2017. — 336 с.: ил.

ISBN 978-5-9775-3758-2

Книга позволяет изучить науку о данных (Data Science) и применить полученные знания на практике. Она написана так, что способствует погружению в Data Science аналитика, фактически не обладающего глубокими знаниями в этой прикладной дисциплине.

В объемах, достаточных для начала работы в области Data Science, книга содержит интенсивный курс языка Python, элементы линейной алгебры, математической статистики, теории вероятностей, методов сбора, очистки, нормализации и обработки данных. Даны основы машинного обучения. Описаны различные математические модели и их реализация по методу k ближайших соседей, наивной байесовской классификации, линейной и логистической регрессии, а также модели на основе деревьев принятия решений, нейронных сетей и кластеризации. Рассказано о работе с рекомендательными системами, описаны приемы обработки естественного языка, методы анализа социальных сетей, основы баз данных, SQL и MapReduce.

*Для аналитиков данных*

УДК 004.6  
ББК 32.81

**Группа подготовки издания:**

Главный редактор	<i>Екатерина Кондукова</i>
Зам. главного редактора	<i>Евгений Рыбаков</i>
Зав. редакцией	<i>Екатерина Капалыгина</i>
Перевод с английского	<i>Андрея Логунова</i>
Редактор	<i>Анна Кузьмина</i>
Компьютерная верстка	<i>Ольги Сергиенко</i>
Корректор	<i>Зинаида Дмитриева</i>
Оформление обложки	<i>Марины Дамбиевой</i>

Authorized translation of the English edition of Data Science from Scratch (978-1-491-90142-7) © 2015 Joel Grus.  
This translation is published and sold by permission of O'Reilly Media, Inc., which owns or controls all rights to publish and sell the same.

Авторизованный перевод английской редакции книги Data Science from Scratch (ISBN: 978-1-491-90142-7)  
© 2015 Joel Grus.

Перевод опубликован и продается с разрешения O'Reilly Media, Inc., собственника всех прав на публикацию и продажу издания.

Подписано в печать 28.04.17.  
Формат 70×100<sup>1/16</sup>. Печать офсетная. Усл. печ. л. 27,09.  
Тираж 1000 экз. Заказ № 4330.  
"БХВ-Петербург", 191036, Санкт-Петербург, Гончарная ул., 20.

ООО "Печатное дело",  
142300, МО, г. Чехов, ул. Полиграфистов, д. 1

ISBN 978-1-491-90142-7 (англ.)  
ISBN 978-5-9775-3758-2 (рус.)

© 2015 Joel Grus  
© Перевод на русский язык "БХВ-Петербург", 2017

---

# Data Science from Scratch

*Joel Grus*

Beijing • Cambridge • Farnham • Sebastopol • Tokyo

**O'REILLY**

---

# Оглавление

<b>Предисловие</b> .....	<b>11</b>
Наука о данных .....	11
С чистого листа .....	12
Условные обозначения, принятые в книге .....	13
Использование примеров кода .....	14
Благодарности .....	15
<b>Комментарий переводчика</b> .....	<b>16</b>
Python 2 и Python 3 .....	16
Установка и удаление дистрибутива Anaconda .....	17
Настройка дистрибутива Anaconda .....	18
Установка инструментальной среды Spyder .....	18
Настройка инструментальной среды Spyder .....	19
Настройка среды Spyder с Python 3 для работы с Python 2 .....	19
Факультативно .....	20
Запуск сервера записных книжек Jupyter .....	20
Установка библиотек Python из whl-файла .....	20
Подготовка среды Python 3 в ОС Ubuntu Linux .....	21
Управление пакетами .deb в Ubuntu Linux .....	21
<b>Об авторе</b> .....	<b>23</b>
<b>Глава 1. Введение</b> .....	<b>25</b>
Господство данных .....	25
Что такое наука о данных? .....	25
Оправдание для выдумки: DataSciencester .....	27
Поиск ключевых звеньев .....	27
Аналитики, которых вы должны знать .....	30
Зарплаты и опыт работы .....	33
Оплата премиум-аккаунтов .....	35
Популярные темы .....	36
Вперед .....	38
<b>Глава 2. Интенсивный курс языка Python</b> .....	<b>39</b>
Основы .....	39
Установка .....	39
Дзен языка Python .....	40
Пробельные символы .....	40
Модули .....	41
Арифметические операции .....	42

Функции.....	43
Строки.....	44
Исключения.....	44
Списки.....	45
Кортежи.....	46
Словари.....	47
Словарь <i>defaultdict</i> .....	48
Словарь <i>Counter</i> .....	49
Множества.....	50
Управляющие конструкции.....	50
Истинность.....	51
Не совсем основы.....	52
Сортировка.....	52
Генераторы последовательностей.....	53
Функции-генераторы и генераторные выражения.....	54
Случайные числа.....	55
Регулярные выражения.....	56
Объектно-ориентированное программирование.....	56
Инструменты функционального программирования.....	58
Функция <i>enumerate</i> .....	59
Функция <i>zip</i> и распаковка аргументов.....	60
Переменные <i>args</i> и <i>kwargs</i> .....	60
Добро пожаловать в DataSciencester!.....	62
Для дальнейшего изучения.....	62
<b>Глава 3. Визуализация данных.....</b>	<b>63</b>
Библиотека <i>matplotlib</i> .....	63
Столбчатые диаграммы.....	65
Линейные графики.....	68
Точечные диаграммы.....	70
Для дальнейшего изучения.....	72
<b>Глава 4. Линейная алгебра.....</b>	<b>73</b>
Векторы.....	73
Матрицы.....	77
Для дальнейшего изучения.....	80
<b>Глава 5. Статистика.....</b>	<b>81</b>
Описание одиночного набора данных.....	81
Показатели центра распределения.....	83
Показатели вариации.....	85
Корреляция.....	87
Парадокс Симпсона.....	90
Некоторые другие ловушки корреляции.....	91
Корреляция и причинная зависимость.....	91
Для дальнейшего изучения.....	92
<b>Глава 6. Теория вероятностей.....</b>	<b>93</b>
Зависимость и независимость.....	93
Условная вероятность.....	94

Теорема Байеса .....	96
Случайные величины.....	97
Непрерывные распределения.....	98
Нормальное распределение .....	100
Центральная предельная теорема.....	103
Для дальнейшего изучения .....	105
<b>Глава 7. Гипотеза и вывод.....</b>	<b>106</b>
Проверка статистических гипотез.....	106
Пример: бросание монеты .....	106
<i>P</i> -значения .....	110
Доверительные интервалы .....	111
Подгонка <i>p</i> -значения .....	112
Пример: проведение <i>A/B</i> -тестирования .....	113
Байесовский статистический вывод.....	115
Для дальнейшего изучения .....	118
<b>Глава 8. Градиентный спуск.....</b>	<b>119</b>
Идея в основе метода градиентного спуска .....	119
Вычисление градиента .....	120
Использование градиента .....	123
Выбор оптимального размера шага .....	124
Собираем все вместе .....	124
Стохастический градиентный спуск .....	126
Для дальнейшего изучения .....	127
<b>Глава 9. Сбор данных .....</b>	<b>129</b>
Объекты <i>stdin</i> и <i>stdout</i> .....	129
Чтение файлов.....	131
Основы работы с текстовыми файлами .....	131
Файлы с разделителями.....	132
Извлечение данных из веб-ресурсов .....	134
Анализ кода HTML .....	134
Пример: книги об анализе данных издательства O'Reilly .....	137
Использование программных интерфейсов .....	141
Формат JSON (и XML).....	141
Использование непроверенного API .....	142
Поиск API .....	144
Пример: использование интерфейсов Twitter API .....	144
Получение учетных данных .....	145
Использование Twython .....	146
Для дальнейшего изучения .....	148
<b>Глава 10. Обработка данных.....</b>	<b>149</b>
Исследование данных.....	149
Исследование одномерных данных.....	149
Двумерные данные .....	151
Многомерные данные.....	153
Очистка и форматирование .....	155
Управление данными .....	157

Шкалирование.....	160
Снижение размерности .....	162
Для дальнейшего изучения .....	168
<b>Глава 11. Машинное обучение.....</b>	<b>169</b>
Моделирование .....	169
Что такое машинное обучение?.....	170
Переобучение и недообучение .....	171
Правильность модели.....	173
Компромисс между смещением и дисперсией.....	176
Извлечение и отбор признаков .....	177
Для дальнейшего изучения .....	178
<b>Глава 12. К ближайших соседей .....</b>	<b>180</b>
Модель.....	180
Пример: предпочтительные языки .....	182
Проблема проклятия размерности .....	186
Для дальнейшего изучения .....	190
<b>Глава 13. Наивный Байес .....</b>	<b>191</b>
Действительно глупый спам-фильтр.....	191
Более продуманный спам-фильтр .....	192
Реализация.....	194
Тестирование модели .....	196
Для дальнейшего изучения .....	198
<b>Глава 14. Простая линейная регрессия .....</b>	<b>199</b>
Модель.....	199
Применение метода градиентного спуска .....	202
Метод максимального правдоподобия .....	203
Для дальнейшего изучения .....	204
<b>Глава 15. Множественная регрессия.....</b>	<b>205</b>
Модель.....	205
Другие допущения модели наименьших квадратов.....	206
Подбор модели.....	207
Интерпретация модели.....	208
Качество подбора модели .....	209
Отступление: бутстрапирование данных.....	209
Стандартные ошибки коэффициентов регрессии .....	211
Регуляризация .....	213
Для дальнейшего изучения .....	215
<b>Глава 16. Логистическая регрессия.....</b>	<b>216</b>
Задача.....	216
Логистическая функция .....	218
Применение модели.....	220
Качество подбора модели .....	221
Метод опорных векторов .....	223
Для дальнейшего изучения .....	225



<b>Глава 17. Деревья принятия решений</b> .....	<b>226</b>
Что такое дерево принятия решений? .....	226
Энтропия .....	228
Энтропия разбиения .....	230
Создание дерева принятия решений .....	231
Собираем все вместе .....	233
Случайные леса .....	236
Для дальнейшего изучения .....	237
<b>Глава 18. Нейронные сети</b> .....	<b>238</b>
Перцептроны .....	238
Нейронные сети прямого распространения .....	240
Метод обратного распространения ошибки .....	243
Пример: преодоление капчи .....	244
Для дальнейшего изучения .....	249
<b>Глава 19. Кластеризация</b> .....	<b>250</b>
Идея .....	250
Модель .....	251
Пример: встречи для специалистов .....	252
Выбор числа $k$ .....	254
Пример: кластеризация цвета .....	256
Восходящий метод иерархической кластеризации .....	257
Для дальнейшего изучения .....	263
<b>Глава 20. Обработка естественного языка</b> .....	<b>264</b>
Облака слов .....	264
$N$ -граммные модели языка .....	266
Грамматики .....	269
Ремарка: метод сэмплирования по Гиббсу .....	271
Тематическое моделирование .....	273
Для дальнейшего изучения .....	278
<b>Глава 21. Анализ социальных сетей</b> .....	<b>279</b>
Центральность по посредничеству .....	279
Центральность собственного вектора .....	284
Умножение матриц .....	284
Центральность .....	287
Направленные графы и рейтинг PageRank .....	288
Для дальнейшего изучения .....	291
<b>Глава 22. Рекомендательные системы</b> .....	<b>292</b>
Неавтоматическое кураторство .....	293
Рекомендация популярных тем .....	293
Коллаборативная фильтрация на основе пользователя .....	294
Коллаборативная фильтрация по схожести предметов .....	297
Для дальнейшего изучения .....	300
<b>Глава 23. Базы данных и SQL</b> .....	<b>301</b>
Операторы <i>CREATE TABLE</i> и <i>INSERT</i> .....	301
Оператор <i>UPDATE</i> .....	303

Оператор <i>DELETE</i> .....	304
Оператор <i>SELECT</i> .....	304
Оператор <i>GROUP BY</i> .....	306
Оператор <i>ORDER BY</i> .....	308
Оператор <i>JOIN</i> .....	309
Подзапросы .....	311
Индексы .....	312
Оптимизация запросов .....	313
Базы данных NoSQL .....	313
Для дальнейшего изучения .....	314
<b>Глава 24. Распределенные вычисления MapReduce</b> .....	<b>315</b>
Пример: подсчет частотности слов .....	315
Почему MapReduce? .....	317
MapReduce в более общей реализации .....	318
Пример: анализ обновлений ленты новостей .....	319
Пример: умножение матриц .....	321
Ремарка: сумматоры .....	322
Для дальнейшего изучения .....	323
<b>Глава 25. Идите и займитесь аналитикой</b> .....	<b>324</b>
Интерактивная оболочка IPython .....	324
Математический аппарат .....	325
Не с чистого листа .....	325
Библиотека NumPy .....	326
Библиотека pandas .....	326
Библиотека scikit-learn .....	326
Визуализация .....	326
Язык программирования R .....	327
Где найти данные? .....	327
Занятия анализом данных .....	328
Новости хакера .....	328
Пожарные машины .....	329
Футболки .....	329
А вы? .....	330
<b>Предметный указатель</b> .....	<b>331</b>

---

# Предисловие

## Наука о данных

Аналитиков данных (data scientists) называют "самой сексуальной профессией XXI века". Очевидно тот, кто так выразился, никогда не бывал в пожарной части. Тем не менее, наука о данных (data science) — это действительно передовая и быстро развивающаяся отрасль знаний, а чтобы отыскать обозревателей рыночных тенденций, которые возбужденно предвещают, что через 10 лет нам потребуются на миллиарды и миллиарды больше аналитиков данных, чем мы имеем на текущий момент, не придется долго рыскать по Интернету.

Но что же это такое — наука о данных? В конце концов нельзя же выпускать специалистов в этой области, если не знаешь, что она собой представляет. Согласно диаграмме Венна<sup>1</sup>, которая довольно известна в этой отрасли, наука о данных находится на пересечении:

- ◆ навыков алгоритмизации и программирования;
- ◆ знаний математики и статистики;
- ◆ профессионального опыта в предметной области.

Собираясь с самого начала написать книгу, охватывающую все три направления, я быстро пришел к выводу, что всестороннее обсуждение профессионального опыта в предметной области потребовало бы десятки тысяч страниц. Тогда я решил сосредоточиться на первых двух. Задача — помочь желающим развить свои навыки алгоритмизации и программирования, которые потребуются для того, чтобы приступить к решению задач в области науки о данных, а также почувствовать себя комфортно с математикой и статистикой, которые находятся в центре этой междисциплинарной практической сферы.

Довольно трудновыполнимая задача для книги. Развивать свои навыки алгоритмизации и программирования лучше всего решая прикладные задачи. Прочтя эту книгу, читатель получит хорошее представление о моих способах алгоритмизации прикладных задач, некоторых инструментах, которыми я пользуюсь в работе, и моем подходе к решению задач, связанных с анализом данных. Но это вовсе не означает, что мои способы, инструменты и подход являются единственно возможными. Надеюсь, что мой опыт вдохновит попробовать пойти собственным путем. Все ис-

---

<sup>1</sup> Джон Венн (1834–1923) — английский логик, предложивший схематичное изображение всех возможных пересечений нескольких (часто — трех) множеств (см. [https://ru.wikipedia.org/wiki/Венн,\\_Джон](https://ru.wikipedia.org/wiki/Венн,_Джон)). — *Прим. пер.*

ходные коды и данные из книги доступны на GitHub (<https://github.com/joelgrus/data-science-from-scratch>); они помогут приступить к работе.

То же самое касается и математики. Ею лучше всего заниматься, решая математические задачи. Данная книга, разумеется, не является руководством по математике, и мы не собираемся, по большей части, заниматься ею. Однако действительность такова, что заниматься анализом данных не получится без *некоторых* представлений о теории вероятностей, математической статистике и линейной алгебре. Это значит, что по мере надобности мы будем уходить с головой в математические равенства, интуицию, аксиомы и мультяшные версии глобальных математических идей. Надеюсь, читатель не побоится погрузиться в них вместе со мной.

По ходу изложения я также надеюсь передать ощущение, что импровизировать с данными — это увлекательное занятие, потому что, как бы это сказать, импровизировать с ними на самом деле увлекательно! В особенности, если сравнивать с такими альтернативами, как подготовка налоговой декларации или добыча угля.

## С чистого листа

Для работы в области науки о данных разработана масса программных библиотек, платформ, модулей и инструментариев, которые эффективно реализуют наиболее (нередко, и наименее) общие алгоритмы и приемы, применяемые в науке о данных. Тот, кто станет аналитиком данных, несомненно, будет досконально знать библиотеку для научных вычислений NumPy, библиотеку для машинного обучения scikit-learn, библиотеку для анализа данных pandas и множество других. Они прекрасно подходят для решения задач, связанных с наукой о данных. Но они также способствуют тому, чтобы начать решать задачи в области науки о данных, фактически не понимая ее.

В этой книге мы начнем двигаться в сторону науки о данных, стартовав с нулевой отметки, а именно займемся разработкой инструментов и реализацией алгоритмов вручную с тем, чтобы лучше понять их. Я вложил немало своего умственного труда в создание ясных, хорошо задокументированных и читаемых реализаций алгоритмов и примеров. В большинстве случаев инструменты, которые мы станем конструировать, будут иметь не практический, а разъясняющий характер. Они хорошо работают с малыми, почти игрушечными, наборами данных, но не справляются с данными "веб-масштаба".

По ходу изложения я буду отсылать читателя к библиотекам, которые подойдут для применения этих методов на более крупных наборах данных. Но мы их не будем тут использовать.

По поводу того, какой язык программирования лучше всего подходит для обучения науке о данных, развернулась здоровая полемика. Многие настаивают на языке статистического программирования R (мы называем таких людей неправильными). Некоторые предлагают Java или Scala. Однако, по моему мнению, Python — идеальный вариант.

Он обладает несколькими особенностями, которые делают его особенно пригодным для изучения и решения задач в области науки о данных:

- ◆ он бесплатный;
- ◆ он относительно прост в написании кода (и в особенности в понимании);
- ◆ он располагает сотнями прикладных библиотек, предназначенных для работы в области науки о данных.

Не решусь назвать Python моим предпочтительным языком программирования. Есть другие языки, которые я нахожу более удобными, продуманными либо просто интересными для программирования. И все же практически всякий раз, когда я начинаю новый проект в области науки о данных, либо, когда мне нужно быстро набросать рабочий прототип, либо продемонстрировать концепции этой практической дисциплины ясным и легким для понимания способом, я всякий раз в итоге использую Python. И поэтому в этой книге используется Python.

Эта книга не предназначена для того, чтобы научить программировать на Python (хотя, я почти уверен, что, прочтя ее, этому можно немного научиться). Тем не менее, я проведу интенсивный курс программирования на Python (ему будет посвящена целая глава), где будут высвечены характерные черты, которые в данном случае приобретают особую важность. Однако если знания, как программировать на Python (или о программировании вообще), отсутствуют, то читателю остается самому дополнить эту книгу чем-то вроде руководства по Python для начинающих.

В последующих частях нашего введения в науку о данных будет принят такой же подход — углубляться в детали там, где они оказываются решающими или показательными. В других ситуациях на читателя возлагается задача домысливать детали самому или заглядывать в Википедию.

За годы работы в отрасли я подготовил некоторое количество специалистов в области науки о данных. Хотя не все из них стали меняющими мир супер-мега-рок-звездами в области анализа данных, тем не менее, я их всех выпустил более подготовленными специалистами, чем они были до этого. И я все больше убеждаюсь в том, что любой, у кого есть некоторые математические способности вкупе с определенным набором навыков в программировании, располагает всем необходимым для решения задач в области науки о данных. Все, чего она требует, — это лишь пытливый ум, готовность усердно трудиться и наличие данной книги. Отсюда и эта книга.

## Условные обозначения, принятые в книге

В книге используются следующие типографические условные обозначения:

- ◆ *курсив* указывает на новые термины;
- ◆ моноширинный шрифт используется для листингов программ, а также внутри абзацев для ссылки на элементы программ, такие как переменные или имена функций, базы данных, типы данных, переменные окружающей среды, операторы и ключевые слова;

- ◆ **жирный моноширинный шрифт** показывает команды либо другой текст, который должен быть напечатан пользователем, а также ключевые слова в коде;
- ◆ **моноширинный шрифт курсивом** показывает текст, который должен быть заменен значениями пользователя либо значениями, определяемыми по контексту.



Данный элемент обозначает подсказку или совет.



Данный элемент обозначает общее замечание.



Данный элемент обозначает предупреждение или предостережение.

## Использование примеров кода

Дополнительный материал (примеры кода, упражнения и пр.) доступен для скачивания по адресу <https://github.com/joelgrus/data-science-from-scratch>.

Эта книга предназначена для того, чтобы помочь вам решить ваши задачи. В целом, если код примеров предлагается вместе с книгой, то вы можете использовать его в своих программах и документации. Вам не нужно связываться с нами с просьбой о разрешении, если вы не воспроизводите значительную часть кода. Например, написание программы, которая использует несколько фрагментов кода из данной книги, официального разрешения не требует.

Адаптированный вариант примеров в виде электронного архива вы можете скачать по ссылке <ftp://ftp.bhv.ru/9785977537582.zip>. Эта ссылка доступна также со страницы книги на сайте [www.bhv.ru](http://www.bhv.ru).

После распаковки архива у вас получится несколько папок:

- ◆ папки `code-python2-original` и `code-python3-original` содержат исходные примеры автора соответственно для Python 2 и Python 3;
- ◆ папка `code-python3-ru` — адаптированные исходные примеры программ;
- ◆ папка "Записные книжки Jupyter" — соответственно записные книжки Jupyter и html-версии книжек;
- ◆ папка Прочее — файлы корпуса текстов SpamAssassin и whl-файл библиотеки BeautifulSoup, которая понадобится в *главе 10*, как образец whl-пакета, и которую впрочем можно скачать и установить самостоятельно (см. разд. "Комментарии переводчика").

## Благодарности

Прежде всего, хотел бы поблагодарить Майка Лукидеса (Mike Loukides) за то, что принял мое предложение по поводу этой книги, и за то, что настоял на том, чтобы я довел ее до разумного объема. Он мог бы легко сказать: "Кто этот человек, вообще, который шлет мне образцы глав, и как заставить его прекратить, наконец?" И я признателен, что он этого не сделал. Хотел бы поблагодарить моего редактора, Мари Богуро (Marie Beaugureau), которая консультировала меня в течение всей процедуры публикации и привела книгу в гораздо более удобоваримый вид, чем если бы этим я занимался сам.

Я бы не написал эту книгу, если бы не изучил науку о данных, и, возможно, я бы не научился ей, если бы не влияние Дэйва Хсу (Dave Hsu), Игоря Татарина (Igor Tatarinov), Джона Раузера (John Rauser) и остальной группы единомышленников из Farecast (это было так давно, что в то время даже не было понятия науки о данных). Хорошие парни в Coursera также заслуживают много добрых слов.

Я так же благодарен моим бета-читателям и рецензентам. Джэй Фандлинг (Jay Fundling) обнаружил тонны ошибок и указал на многие нечеткие объяснения, и в итоге книга оказалась намного лучше и корректней, благодаря ему. Дэбаши Гош (Debashis Ghosh) героически провела санитарную проверку моих статистических выкладок. Эндрю Массельман (Andrew Musselman) предложил смягчить первоначальный аспект книги, касающийся утверждения, что "люди, предпочитающие язык R вместо Python, есть моральные извращенцы", что, думаю, в итоге оказалось неплохим советом. Трэй Кози (Trey Causey), Райан Мэттью Балфанц (Ryan Matthew Balfanz), Лорис Муларони (Loris Mularoni), Нуриа Пуджол (Núria Pujol), Роб Джефферсон (Rob Jefferson), Мэри Пэт Кэмпбелл (Mary Pat Campbell), Зак Джири (Zach Geary) и Венди Грас (Wendy Grus) также высказали неоценимые замечания. За любые оставшиеся ошибки, конечно же, отвечаю только я.

Я многим обязан сообществу с хештегом **#datascience** в Twitter за то, что продемонстрировали мне массу новых концепций, познакомили меня со множеством замечательных людей и заставили почувствовать себя двоечником, да так, что я ушел и написал книгу в качестве противовеса. Особые благодарности снова Трэйу Кози (Trey Causey) за то, что нечаянно упомянул, чтобы я включил главу про линейную алгебру, и Шин Дж. Тэйлор (Sean J. Taylor) за то, что неумышленно указала на пару больших несоответствий в главе "Работа с данными".

Но больше всего я задолжал огромную благодарность Ганге (Ganga) и Мэдлин (Madeline). Единственная вещь, которая труднее написания книги, — это жить рядом с тем, кто пишет книгу, и я бы не смог ее закончить без их поддержки.

---

# Комментарий переводчика

Приведенные в книге примеры были протестированы в инструментальной среде для научных вычислений для языка Python — Spyder (Scientific PYthon Development EnviRonment) — версии 3.0.0 на основе Python 3.5.2 для Windows. Spyder — это простая, легковесная и, главное, бесплатная, интерактивная среда разработки на Python, которая предлагает функции, похожие на среду разработки в MATLAB. Она также предоставляет написанные полностью на Python и готовые к использованию виджеты PyQt4 и PySide, редактор исходного кода с подсветкой синтаксиса и функциями самоанализа/анализа кода, редактор массивов данных NumPy, редактор словарей, консоли Python и IPython и многое другое. Большинство графиков и диаграмм были воссозданы именно в этой инструментальной среде, хотя стоит признать, есть и другие, например PyCharm и Eclipse PyDev.

В связи с тем, что используемые в примерах имена переменных и функций являются логическим продолжением излагаемого материала, они переведены вместе с авторскими комментариями и дополнены кратким описанием. Кроме того, в силу различий в наименовании терминов как в отечественной, так и зарубежной литературе, для некоторых из них приведены соответствующие эквивалентные варианты наименований или пояснения. Некоторые термины, которые в тексте остались без объяснения, кратко определены в сносках.

В целом в книге принят сбалансированный подход к теории и программированию. Разумеется, можно найти более подробное изложение аспектов науки о данных и на основе других языков программирования (Clojure, Julia, Haskell, MATLAB, R или Scala). Однако если необходимо познакомиться с этой тематикой на основе языка Python, то, скорее всего, начинать следует именно с этой книги.

Книга может быть интересной широкому кругу специалистов, в том числе в области машинного обучения, начинающим аналитикам данных, преподавателям, студентам, а также всем, кто интересуется программированием.

## Python 2 и Python 3

Сегодня используются две главные разновидности Python: Python версии 2.x существует уже в течение многих лет и все еще широко применяется, в то время как Python версии 3.x, которая не является обратно несовместимой с Python 2, становится все более популярной, т. к. эта версия Python взята за основу для дальнейшего развития.

Один из главных факторов, сдерживающих повсеместное принятие Python 3, — недостаточная поддержка сторонних библиотек, в частности связанных с разработкой



геоприложений или визуальных интерфейсов. Но мир не стоит на месте, а вместе с ним и Python 3, и все крупнейшие библиотеки, которые указаны в этой книге, теперь в равной степени можно выполнять, используя Python 3. Все примеры программного кода в этой книге преобразованы таким образом, чтобы использовать синтаксис Python 3. К счастью, различия в синтаксисе между Python 2 и Python 3 предельно простые и потребуют незначительных изменений. В простых примерах часто единственное отличие состоит в том, что в версии Python 3 после оператора `print` требуются круглые скобки; все остальные отличия задокументированы в сносках.

## Установка и удаление дистрибутива Anaconda

Anaconda — это полностью свободный дистрибутив Python (предназначенный в том числе для коммерческого использования и повторного распространения). Он содержит более 400 самых популярных библиотек Python для вычислений в области естественных наук, математики, инженерии и анализа данных.

Установка дистрибутива в Windows выполняется стандартным образом после загрузки с сайта бинарного файла дистрибутива. Папка с установленными приложениями находится в разделе приложений меню **Пуск**. Откройте папку меню **Пуск**, выберите **Все приложения** и нажмите на папке **Anaconda3 (Пуск | Все приложения | Anaconda3)**. Папка содержит ряд приложений, таких как навигатор Anaconda Navigator, инструментальная среда Spyder и сервер записных книжек Jupyter. Исполнимые файлы находятся в папке `C:\Users\имя_пользователя\Anaconda3\`.

В Linux сначала необходимо скачать установщик — скриптовый файл с расширением `.sh` для оболочки Bash ([https://www.continuum.io/downloads#\\_unix](https://www.continuum.io/downloads#_unix)). На примере дистрибутива Anaconda версии 3 (для Python 3.5.2) команда установки выглядит так:

```
bash Anaconda3-4.1.1-Linux-x86_64.sh
```

Отметим, что инсталляция вступит в силу после того, как вы закроете и заново откроете окно терминала.

Для обновления дистрибутива Anaconda следует набрать в терминале следующую команду:

```
conda update conda
```

Удаление дистрибутива Anaconda:

```
rm -rf ~/anaconda
```

Более подробная информация по деинсталляции Anaconda содержится по указанной выше ссылке.

Чтобы удостовериться, что дистрибутив Anaconda установлен успешно, воспользуйтесь следующей командой:

```
conda --version
```

В результате будет выведен номер установленной версии.

## Настройка дистрибутива Anaconda

В случае если вы в основном работаете в среде Python 3, но иногда возникает необходимость переключиться в среду Python 2 и использовать ее для работы с библиотеками, которые предназначены только для Python 2, то Anaconda предлагает функционал создания и активации новой среды, куда можно установить нужную версию языка. В табл. 1 приведено несколько команд, которые можно выполнить прямо в Spyder во встроенном окне командной строки (Инструменты | Открыть командную строку).

Таблица 1

Команда	Описание
<code>conda create -n py27 python=2.7.9 anaconda</code>	Установить другую версию Python (2.7.9) в новую среду с именем py27 (имя может быть любым)
<code>conda info --envs</code>	Проверить, что среда с именем py27 установлена (команда выводит список всех сред, при этом активная среда выделяется знаком *)
<code>source activate py27 (Linux, OS X), activate py27 (Windows)</code>	Переключиться в среду py27 с другой версией Python (команда <code>activate</code> добавляет в начало строки путь к среде py27)
<code>deactivate</code>	Деактивировать текущую среду и вернуться к среде по умолчанию
<code>python --version</code>	Проверить, что среда использует нужную версию Python
<code>conda remove -n py27 --all</code>	Удалить среду py27 (после выполнения команды выполнить <code>conda clean --lock</code> )
<code>conda install --name py27 scyru</code>	Установить библиотеку scyru в среду py27
<code>conda remove --name py27 scyru</code>	Удалить библиотеку scyru в среде py27
<code>conda clean --lock</code>	Очистить блокировку, если произошел сбой при установке среды (в Windows иногда сперва требуется удалить conda в диспетчере задач)

## Установка инструментальной среды Spyder

Для пользователей Windows хорошая новость заключается в том, что инструментальная среда программирования Spyder уже включена в состав дистрибутива Anaconda, и исполнимый файл находится в папке `C:\Users\[имя_пользователя]\Anaconda3\Scripts\spyder.exe`.

Чтобы установить среду Spyder в Ubuntu Linux, используя официальный менеджер пакетов, нужна всего одна команда (тройка соответствует Python версии 3):

```
sudo apt-get install spyder
```

Чтобы установить с использованием менеджера пакетов `pip`:

```
sudo apt-get install python3-qt4 python3-sphinx
sudo pip3 install spyder
```

И чтобы обновить:

```
sudo pip3 install -U spyder
```

Для выполнения команд менеджера пакетов `pip` следует открыть консольное окно (**Инструменты | Открыть командную строку**) и в открывшемся в правой части интерфейса пользователя окне набрать нужную команду, например, `pip3 install jupyter`.

## Настройка инструментальной среды Spyder

При инсталляции дистрибутива Anaconda3 по умолчанию базовым является интерпретатор Python 3 (в данном случае версии 3.5.2). В случае если нужно на время переключиться в режим работы в интерпретаторе Python версии 2 (в данном случае версии 2.7.9), существует два варианта: дополнительно установить дистрибутив Anaconda для Python 2 либо выполнить небольшую настройку инструментальной среды Spyder, по умолчанию работающей на основе Python 3.

## Настройка среды Spyder с Python 3 для работы с Python 2

Для такой настройки нужно в основном меню выбрать **Инструменты** и затем **Параметры**. В открывшемся окне **Параметры** выбрать **Интерпретатор Python**. В разделе **Интерпретатор Python** с двумя переключателями установить переключатель **Использовать следующий интерпретатор Python** и затем нажать кнопку напротив текстового поля, чтобы выбрать путь к интерпретатору Python 2.7.9, который находится в папке `C:\Users\имя_пользователя\Anaconda3\envs\py27`, где `py27` — это имя среды, созданной вами для Python 2.7.9 (см. разд. "Настройка дистрибутива Anaconda" ранее). Если открыть новую консоль (**Консоли | Открыть консоль Python**), то в строке приветствия будет указана нужная версия Python.

Теперь в этой консоли можно набирать команды и вставлять целые программы, однако запускать программы из рабочего окна редактора Spyder не получится. Для сценария, работающего с Python 2, требуется еще одна и последняя настройка. В основном меню нужно выбрать **Запуск**, затем **Настроить** и в разделе **Консоль** выбрать второй переключатель **Выполнить во внешнем системном терминале** (т. е. в новой специально выделенной консоли). Теперь этот конкретный сценарий будет выполняться в консоли с Python 2.

Закончив работу с интерпретатором Python 2.7.9, не забудьте вернуться к исходному интерпретатору Python 3.5.2. Для этого нужно в разделе **Интерпретатор Python** просто установить переключатель **По умолчанию** (тот же, что и для Spyder).

Напомним, что в Windows месторасположение базового интерпретатора следующее: `C:\Users\[имя_пользователя]\Anaconda3\python.exe`. Разумеется, такой режим работы вносит некоторые неудобства, но по крайней мере он не такой громоздкий, как установка еще одной версии Anaconda 4.1.1, но уже с Python 2.7.9 в качестве базового интерпретатора.

## Факультативно

### Запуск сервера записных книжек Jupyter

В Windows и Ubuntu Linux локальный сервер записных книжек запускается из командной строки двумя способами. Первый: откройте окно командной оболочки (cmd) в Windows или окно терминала в Ubuntu Linux и просто наберите:

```
jupyter notebook
```

Второй вариант:

```
ipython notebook
```

Теперь интерактивная вычислительная среда IPython — это составная часть интегрированной среды Jupyter, и поэтому второй вариант запуска не рекомендован к применению и будет в будущем удален, к тому же в нем используется формат записной книжки предыдущей версии 3.0+ (текущая 5.0+).

### Установка библиотек Python из whl-файла

Библиотеки для Python можно разрабатывать не только на чистом Python. Довольно часто библиотеки пишутся на C (динамические библиотеки) и для них пишется обертка Python или же библиотека пишется на Python, но для оптимизации узких мест часть кода пишется на C. Такие библиотеки получаются очень быстрыми, однако библиотеки с вкраплениями кода на C программисту на Python тяжелее установить ввиду банального отсутствия соответствующих знаний либо необходимых компонентов и настроек в рабочей среде (в особенности в Windows). Для решения описанных проблем разработан специальный формат (файлы с расширением whl) для распространения библиотек, который содержит заранее скомпилированную версию библиотеки со всеми ее зависимостями. Формат WHL поддерживается всеми основными платформами (Mac OS X, Linux, Windows).

Установка производится с помощью менеджера пакетов pip3 (тройка соответствует Python версии 3). В отличие от обычной установки командой `pip3 install <имя_пакета>`, вместо имени пакета указывается путь к whl-файлу: `pip3 install <путь_к_whl_файлу>`. Например,

```
pip3 install C:\temp\networkx-1.11-py2.py3-none-any.whl
```

Либо откройте окно командной строки и при помощи команды `cd` перейдите в каталог, где размещен ваш whl-файл. В Anaconda по умолчанию подобные файлы нахо-

дятся в каталоге Scripts. Просто скопируйте туда ваш whl-файл. В этих случаях полный путь указывать не понадобится. Например,

```
pip3 install networkx-1.11-py2.py3-none-any.whl
```

При выборе пакета важно, чтобы разрядность устанавливаемой библиотеки и разрядность интерпретатора совпадали. Пользователи Windows могут брать whl-файлы с сайта <http://www.lfd.uci.edu/~gohlke/pythonlibs/>. Библиотеки там постоянно обновляются, и в архиве содержатся все, какие только могут понадобиться.

## Подготовка среды Python 3 в ОС Ubuntu Linux

Таблица 2

Действие	Команда
Обновление ОС Ubuntu/Debian	<code>sudo apt-get update &amp;&amp; sudo apt-get upgrade &amp;&amp; sudo apt-get dist-upgrade &amp;&amp; sudo apt-get autoremove</code>
Интерпретатор языка Python 3 и менеджер пакетов	<code>sudo apt-get install python3 python3-pip</code>
Основные научные библиотеки	<code>sudo apt-get install python3-numpy python3-matplotlib python3-scipy python3-pandas python3-simpy</code>
Записные книжки Jupyter	<code>sudo pip3 install jupyter</code>
OpenGL	<code>sudo apt-get install python3-opengl</code>
Разработка графического интерфейса пользователя	<code>sudo apt-get install python3-pyqt5 python3-pyqt5.qtopenGL python3-pyqt5.qtquick</code>
Хранение данных	<code>sudo apt-get install python3-h5py</code>
Компьютерное зрение	<code>sudo apt-get install python3-skimage</code> <code>sudo apt-get install libatlas-dev libatlas3gf-base &amp;&amp; sudo pip3 install scikit-learn</code>
Интегрированная среда программирования (IDE) Spyder3	<code>sudo apt-get install spyder3</code>
Автозавершение кода в IDE	<code>sudo pip3 install --user rope_py3k</code>

## Управление пакетами .deb в Ubuntu Linux

Установка пакетов .deb в терминале (пакет должен находиться под корневым каталогом home)

```
sudo dpkg -i [пакет].deb
```

или

```
sudo dpkg --install [пакет].deb
```

**Удаление пакетов .deb**

```
sudo dpkg -r [пакет].deb
```

**Удаление вместе с конфигурационными файлами:**

```
sudo dpkg -P [пакет].deb
```

---

## Об авторе

**Джоэл Грас** работает инженером-программистом в компании Google. До этого он занимался аналитической работой в нескольких стартапах. Он живет в Сиэтле, где регулярно посещает неформальные встречи специалистов в области науки о данных. Он редко пишет в свой блог [joelgrus.com](http://joelgrus.com) и всегда доступен в Twitter по хештегу [@joelgrus](https://twitter.com/joelgrus).





# Введение

— Данные! Где данные? — раздраженно восклицал он. —  
Когда под рукой нет глины, из чего лепить кирпичи?

*Артур Конан Дойль<sup>1</sup>*

## Господство данных

Мы живем в мире, страдающем от переизбытка данных. Веб-сайты отслеживают любое нажатие любого пользователя. Смартфоны накапливают сведения о вашем местоположении и скорости в ежедневном и ежесекундном режиме. "Оцифрованные" селферы носят шагомеры на стероидах, которые не переставая записывают их сердечные ритмы, особенности движения, схемы питания и сна. Умные авто собирают сведения о манерах вождения своих владельцев, умные дома — об образе жизни своих обитателей, а умные маркетологи — о наших покупательских привычках. Сам Интернет представляет собой огромный граф знаний, который, среди всего прочего, содержит обширную гипертекстовую энциклопедию, специализированные базы данных о фильмах, музыке, спортивных результатах, игровых автоматах, мемах<sup>2</sup> и коктейлях... и слишком много статистических отчетов (причем некоторые почти соответствуют действительности!) от слишком большого числа государственных исполнительных органов, и все это для того, чтобы вы объяли необъятное.

В этих данных кроются ответы на бесчисленные вопросы, которые никто даже не думает задавать. Эта книга научит вас, как их находить.

## Что такое наука о данных?

Существует шутка, что аналитик данных — это тот, кто знает статистику лучше, чем специалист в области информатики, а информатику — лучше, чем специалист в области статистики. Не утверждаю, что это хорошая шутка, но на самом деле, некоторые аналитики данных действительно являются специалистами в области математической статистики, в то время как другие почти неотличимы от инженеров программного обеспечения. Некоторые являются экспертами в области машинного

---

<sup>1</sup> Реплика Шерлока Холмса из рассказа Артура Конан Дойля (1859–1930) "Медные буки". — *Прим. пер.*

<sup>2</sup> Мем (англ. *meme*) — единица культурной информации: любая идея, символ, манера или образ действия, осознанно или неосознанно передаваемые от человека к человеку посредством речи, письма, видео, ритуалов, жестов и т. д. (см. <https://ru.wikipedia.org/wiki/Мем>). — *Ред.*

обучения, в то время как другие не смогли бы машинно обучиться, чтобы найти выход из детского сада. Некоторые имеют ученые степени доктора наук с впечатляющей историей публикаций, в то время как другие никогда не читали академических статей (хотя, им должно быть стыдно). Короче говоря, в значительной мере неважно, как определять понятие науки о данных, потому что всегда можно найти практикующих аналитиков данных, для которых это определение будет всецело и абсолютно неверным<sup>3</sup>.

Тем не менее, этот факт не остановит нас от попыток. Мы скажем, что аналитик данных — это тот, кто извлекает ценные наблюдения из запутанных данных. В наши дни мир переполнен людьми, которые пытаются превратить данные в ценные наблюдения.

Например, сайт знакомств OkCupid просит своих членов ответить на тысячи вопросов, чтобы отыскать наиболее подходящего для них партнера. Но он также анализирует эти результаты, чтобы вычислить виды безобидных вопросов, с которыми вы можете обратиться, чтобы узнать, насколько высока вероятность близости после первого же свидания.

Компания Facebook просит вас указывать свой родной город и нынешнее местоположение, якобы чтобы облегчить вашим друзьям находить вас и связываться с вами. Но она также анализирует эти местоположения, чтобы определить схемы глобальной миграции и места проживания фанатов различных футбольных команд.

Крупный оператор розничной торговли Target отслеживает покупки и взаимодействия онлайн и в магазине. Он использует данные, чтобы строить прогнозные модели в отношении того, какие клиентки беременны, чтобы лучше продавать им товары, предназначенные для младенцев.

В 2012 г. избирательный штаб Барака Обамы нанял десятки аналитиков данных, которые всюду копали и экспериментировали, чтобы определить избирателей, которым требовалось дополнительное внимание, при этом подбирая оптимальные обращения и программы по привлечению финансовых ресурсов, которые направлялись в адрес конкретных получателей, и сосредотачивая усилия по выводу соперника из предвыборной гонки там, где эти усилия могли быть наиболее успешными. Существует общее мнение, что эти усилия сыграли важную роль в переизбрании президента, вследствие чего совершенно очевидно, что будущие политические кампании будут все более и более управляемыми данными, ведя к бесконечному наращиванию усилий в области науки о данных и методов сбора данных.

И прежде чем вы почувствуете пресыщение, скажем еще пару слов: некоторые аналитики данных время от времени используют свои навыки во благо, чтобы сделать

---

<sup>3</sup> *Наука о данных* — это практическая дисциплина, которая занимается изучением методов обобщаемого извлечения знаний из данных. Она состоит из различных составляющих и основывается на методах и теориях из многих областей знаний, включая обработку сигналов, математику, вероятностные модели, машинное и статистическое обучение, программирование, технологии данных, распознавание образов, теорию обучения, визуальный анализ, моделирование неопределенности, организацию хранилищ данных, а также высокоэффективные вычисления с целью извлечения смысла из данных и создания продуктов обработки данных. — *Прим. пер.*

правительство более эффективным, помочь бездомным и усовершенствовать здравоохранение. И конечно же вы не нанесете вреда своей карьере, если вам нравится заниматься поисками наилучшего способа, как заставить людей щелкать на рекламных баннерах.

## Оправдание для выдумки: DataSciencester

Поздравляем! Вас только что приняли на работу, чтобы вы возглавили усилия в области науки о данных в DataSciencester, *уникальной* социальной сети для аналитиков данных.

Предназначенная исключительно для аналитиков данных, тем не менее, DataSciencester еще ни разу не вкладывалась в развитие собственной практической деятельности в этой области (справедливости ради, она ни разу по-настоящему не вкладывалась даже в развитие своего продукта). Теперь это будет вашей работой! На протяжении всей книги мы будем изучать понятия этой практической дисциплины путем решения задач, с которыми вы будете сталкиваться на работе. Мы будем обращаться к данным, иногда поступающим прямо от пользователей, иногда полученным в результате их взаимодействий с сайтом социальной сети, а иногда даже взятыми из экспериментов, которые будем планировать сами.

И поскольку в DataSciencester царит дух новизны, то мы займемся конструированием собственных инструментов с чистого листа. В результате у вас будет достаточно твердое понимание основ науки о данных, и вы будете готовы применить свои навыки в любой компании или же в любых других задачах, которые окажутся для вас интересными.

Добро пожаловать на борт и удачи! (Вам разрешено носить джинсы по пятницам, и туалет — по коридору направо.)

## Поиск ключевых звеньев

Итак, настал ваш первый рабочий день в DataSciencester, и директор по развитию сети полон вопросов о ее пользователях. До сего дня ему не к кому было обратиться, поэтому он очень воодушевлен вашим прибытием.

В частности, он хочет, чтобы вы установили, кто среди специалистов является "ключевым звеном". Для этих целей он передает вам "снимок" всей социальной сети. (В реальной жизни обычно никто не торопится вручать вам требующиеся данные. Глава 9 посвящена способам сбора данных.)

Как выглядит этот "снимок" данных? Он состоит из списка пользователей `users`, где каждый пользователь представлен ассоциативным списком, или словарем, `dict`, состоящим из идентификатора `id` (т. е. числа) пользователя и его или ее имени `name` (которое по одному из необычайных космических совпадений рифмуется с идентификатором `id`):

```
users = [  
    { "id": 0, "name": "Hero" },  
    { "id": 1, "name": "Dunn" },
```

```
{ "id": 2, "name": "Sue" },
{ "id": 3, "name": "Chi" },
{ "id": 4, "name": "Thor" },
{ "id": 5, "name": "Clive" },
{ "id": 6, "name": "Hicks" },
{ "id": 7, "name": "Devin" },
{ "id": 8, "name": "Kate" },
{ "id": 9, "name": "Klein" }
```

]

Кроме того, он передает вам данные о "дружеских отношениях" friendships в виде списка кортежей, состоящих из пар идентификаторов ID пользователей:

```
friendships = [(0, 1), (0, 2), (1, 2), (1, 3), (2, 3), (3, 4),
               (4, 5), (5, 6), (5, 7), (6, 8), (7, 8), (8, 9)]
```

Например, кортеж (0, 1) означает, что аналитик с id 0 (Hero) и аналитик с id 1 (Dunn) являются друзьями. Сеть представлена на рис. 1.1.

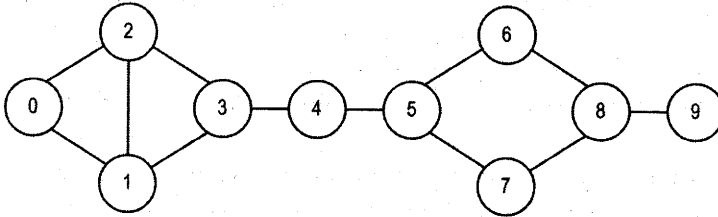


Рис. 1.1. Сеть DataSciencester

Поскольку пользователи представлены словарями dict, то пополнение информации о них становится простым делом.



Не следует прямо сейчас погружаться в подробности реализации кода. В главе 2 будет предложен интенсивный курс языка Python. Пока же следует всего лишь попытаться получить общее представление о том, что мы делаем.

Например, может понадобиться добавить каждому пользователю список друзей. Для этого сначала назначим новому свойству friends каждого пользователя пустой список:

```
# свойство friends содержит друзей для пользователя user
for user in users:
    user["friends"] = []
```

А затем заполним эти списки данными из списка кортежей friendships:

```
for i, j in friendships:
    # это работает, потому что users[i] - это пользователь,
    # чей id равен i
    users[i]["friends"].append(users[j]) # добавить j как друга для i
    users[j]["friends"].append(users[i]) # добавить i как друга для j
```

После того, как в словаре `dict` каждого пользователя размещен список его друзей, получившийся граф можно легко опросить, например, по поводу среднего числа связей.

Сперва находим суммарное число связей, сложив длины всех списков друзей `friends`:

```
# число друзей
def number_of_friends(user):
    """Сколько друзей есть у пользователя user?"""
    return len(user["friends"]) # длина списка id друзей

total_connections = sum(number_of_friends(user) # общее число связей
                        for user in users)     # 24
```

А затем просто делим сумму на число пользователей:

```
from __future__ import division # в Python 2 целочисленное деление хромает
num_users = len(users)          # длина списка пользователей
avg_connections = total_connections / num_users # среднее число связей =
                                         # 2.4
```

Также легко находим наиболее связанных людей, т. е. лиц, имеющих наибольшее число друзей.

Поскольку пользователей не слишком много, их можно упорядочить по убыванию числа друзей:

```
# число друзей для каждого id пользователя
# создать список в формате (id пользователя, число друзей)
num_friends_by_id = [(user["id"], number_of_friends(user))
                    for user in users]

sorted(num_friends_by_id,          # упорядочить его по полю
       key=lambda (user_id, num_friends): num_friends, # num_friends
       reverse=True)              # в убывающем порядке

# каждая пара состоит из id и числа друзей (user_id, num_friends)
# [(1, 3), (2, 3), (3, 3), (5, 3), (8, 3),
#  (0, 2), (4, 2), (6, 2), (7, 2), (9, 1)]
```

Всю проделанную работу можно рассматривать, как способ определить лиц, которые так или иначе занимают в сети центральное место. На самом деле, то, что мы вычислили, представляет собой метрику связи в социальном графе под названием *центральность по степени узлов* (degree centrality) (рис. 1.2).

Достоинство этой метрики заключается в простоте вычислений, однако она не всегда дает нужные или ожидаемые результаты. Например, в `DataSciencester` пользователь `Thor` (id 4) имеет всего две связи, тогда как `Dunn` (id 1) — три. Несмотря на это, схема сети показывает, что `Thor` находится ближе к центру. В *главе 21* мы займемся

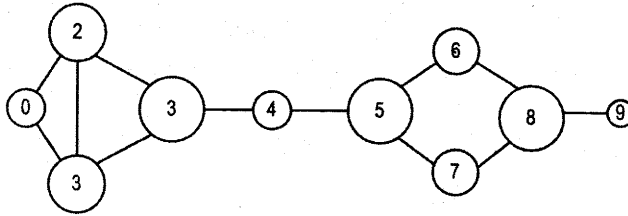


Рис. 1.2. Сеть DataSciencester, измеренная метрикой степени узлов

изучением социальных сетей более подробно и рассмотрим более сложные представления о центральности, которые могут лучше или хуже соответствовать интуиции.

## Аналитики, которых вы должны знать

Вы еще не успели заполнить все страницы формуляра найма новых сотрудников, как директор по поборатимским связям подходит к вашему столу. Она хочет простимулировать рост связей среди членов сети и просит вас разработать систему рекомендации новых друзей "Аналитики, которых Вы должны знать".

Ваша первая реакция — предположить, что пользователь может знать друзей своих друзей. Их легко вычислить: надо лишь просмотреть друзей каждого пользователя и собрать результаты:

```
# список id друзей пользователя user (плохой вариант)
def friends_of_friend_ids_bad(user):
    # "foaf" означает "друг конкретного друга"
    return [foaf["id"]
            for friend in user["friends"] # для каждого из друзей
            # пользователя
            for foaf in friend["friends"]] # получить всех ЕГО друзей
```

Если эту функцию вызвать с аргументом `users[0]` (Неро), она покажет:

```
[0, 2, 3, 0, 1, 3]
```

Список содержит пользователя 0 (два раза), т. к. пользователь Неро на самом деле дружит с обоими своими друзьями; содержит пользователей 1 и 2, хотя оба эти пользователя уже дружат с Неро; и дважды содержит пользователя 3, поскольку Chi достижима через двух разных друзей:

```
print([friend["id"] for friend in users[0]["friends"]]) # [1, 2]
print([friend["id"] for friend in users[1]["friends"]]) # [0, 2, 3]
print([friend["id"] for friend in users[2]["friends"]]) # [0, 1, 3]
```

Информация о том, что некто может стать другом вашего друга несколькими путями, заслуживает внимания, поэтому, напротив, стоит добавить счетчик взаимных друзей, а для этих целей точно потребуется вспомогательная функция, которая будет исключать тех лиц, уже известных пользователю:

```

from collections import Counter # словарь Counter не загружается по умолчанию

# не тот же самый
def not_the_same(user, other_user):
    """два пользователя не одинаковые, если их ключи имеют разные id"""
    return user["id"] != other_user["id"]

# не друзья
def not_friends(user, other_user):
    """other_user - не друг, если он не принадлежит user["friends"], т. е.
    если он not_the_same (не тот же что и все люди в user["friends"])"""
    return all(not_the_same(friend, other_user)
               for friend in user["friends"])

# список id друзей пользователя user
def friends_of_friend_ids(user):
    return Counter(foaf["id"]
                   for friend in user["friends"] # для каждого моего друга
                   for foaf in friend["friends"] # подсчитать ИХ друзей,
                   if not_the_same(user, foaf) # которые не являются мной
                   and not_friends(user, foaf)) # и не мои друзья

print(friends_of_friend_ids(users[3])) # Counter({0: 2, 5: 1})

```

Такая реализация безошибочно сообщает Chi (id 3), что у нее с Nero (id 0) есть двое взаимных друзей, а с Clive (id 5) — всего один такой друг.

Помимо этого, вы понимаете, что, как аналитику данных, вам было бы интересно встречаться с пользователями, которые имеют одинаковую сферу интересов (кстати, это хороший пример аспекта "профессионального опыта в предметной области" науки о данных). После того, как вы поспрашивали вокруг, вам удалось получить на руки данные о сферах интересов в виде списка пар (user\_id, interest), состоящих из id пользователя и интересующей его темы:

```

# интересующие темы
interests = [
    (0, "Hadoop"), (0, "Big Data"), (0, "HBase"), (0, "Java"),
    (0, "Spark"), (0, "Storm"), (0, "Cassandra"),
    (1, "NoSQL"), (1, "MongoDB"), (1, "Cassandra"), (1, "HBase"),
    (1, "Postgres"), (2, "Python"), (2, "scikit-learn"), (2, "scipy"),
    (2, "numpy"), (2, "statsmodels"), (2, "pandas"), (3, "R"), (3, "Python"),
    (3, "statistics"), (3, "regression"), (3, "probability"),
    (4, "machine learning"), (4, "regression"), (4, "decision trees"),
    (4, "libsvm"), (5, "Python"), (5, "R"), (5, "Java"), (5, "C++"),
    (5, "Haskell"), (5, "programming languages"), (6, "statistics"),
    (6, "probability"), (6, "mathematics"), (6, "theory"),
    (7, "machine learning"), (7, "scikit-learn"), (7, "Mahout"),
    (7, "neural networks"), (8, "neural networks"), (8, "deep learning"),

```

```
(8, "Big Data"), (8, "artificial intelligence"), (9, "Hadoop"),
(9, "Java"), (9, "MapReduce"), (9, "Big Data")
```

```
]
```

Например, у Thor (id 4) нет общих друзей с Devin (id 7), но они оба заинтересованы в машинном обучении.

Функция, которая находит пользователей, интересующихся определенной темой, элементарна:

```
# аналитики, которым нравится целевая тема target_interest
def data_scientists_who_like(target_interest):
    return [user_id
            for user_id, user_interest in interests
            if user_interest == target_interest]
```

Все работает, однако функция в такой реализации просматривает весь список тем при каждом запросе. Если пользователей и тем много (либо планируется выполнять много запросов), то лучше создать индексный список пользователей, сгруппированный по теме:

```
from collections import defaultdict
```

```
# id пользователей по значению темы
# ключи - это интересующие темы,
# значения - это списки из id пользователей, интересующихся этой темой
user_ids_by_interest = defaultdict(list)
```

```
for user_id, interest in interests:
    user_ids_by_interest[interest].append(user_id)
```

И еще индексный список тем, сгруппированный по пользователям:

```
# идентификаторы тем по идентификатору пользователя
# ключи - это id пользователей, значения - списки тем для конкретного id
interests_by_user_id = defaultdict(list)
for user_id, interest in interests:
    interests_by_user_id[user_id].append(interest)
```

Теперь легко найти лицо, у кого с конкретным пользователем больше всего общих интересов. Для этого нужно:

1. Выполнить обход интересующих пользователя тем.
2. По каждой теме осуществить обход других пользователей, интересующихся той же самой темой.
3. Подсчитать, сколько раз встретятся другие пользователи.

```
# наиболее общие интересующие темы с пользователем user
def most_common_interests_with(user):
    return Counter(interested_user_id
                  for interest in interests_by_user_id[user["id"]])
```



```
for interested_user_id in user_ids_by_interest[interest]
if interested_user_id != user["id"])
```

Эту функцию можно было бы затем применить в усовершенствованной версии системы рекомендации новых друзей "Аналитики, которых Вы должны знать", основанной на сочетании взаимных друзей и общих тем. Эти виды приложений мы исследуем в *главе 22*.

## Зарплаты и опыт работы

В тот самый момент, когда вы собираетесь пойти пообедать, директор по связям с общественностью обращается к вам с вопросом, можете ли вы предоставить для сайта какие-нибудь примечательные факты об уровне зарплат аналитиков данных. Разумеется, данные о зарплатах имеют конфиденциальный характер, тем не менее, ему удалось снабдить вас анонимным набором данных, содержащим зарплату каждого пользователя (в долларах) и его стаж работы в качестве аналитика данных (в годах):

```
# зарплаты и стаж
salaries_and_tenures = [(83000, 8.7), (88000, 8.1),
                        (48000, 0.7), (76000, 6),
                        (69000, 6.5), (76000, 7.5),
                        (60000, 2.5), (83000, 10),
                        (48000, 1.9), (63000, 4.2)]
```

Первым делом вы выводите данные на диаграмму (в *главе 3* мы рассмотрим, как это делать). Результаты можно увидеть на рис. 1.3.

Совершенно очевидно, что лица с более продолжительным опытом, как правило, зарабатывают больше. Но как это превратить в примечательный факт? Ваша первая попытка — взять среднюю арифметическую зарплату по каждому стажу:

```
# зарплата в зависимости от стажа
# ключи - это годы, значения - это списки зарплат для каждого стажа
salary_by_tenure = defaultdict(list)

for salary, tenure in salaries_and_tenures:
    salary_by_tenure[tenure].append(salary)

# средняя зарплата в зависимости от стажа
# ключи - это годы, каждое значение - это средняя зарплата по этому стажу
average_salary_by_tenure = {
    tenure : sum(salaries) / len(salaries)
    for tenure, salaries in salary_by_tenure.items()
}
```

Но, как оказалось, это не несет какой-то практической значимости, т. к. у всех пользователей разный стаж, и, значит, мы просто сообщаем о зарплатах отдельных пользователей:

```
{0.7: 48000.0,
 1.9: 48000.0,
 2.5: 60000.0,
 4.2: 63000.0,
 6: 76000.0,
 6.5: 69000.0,
 7.5: 76000.0,
 8.1: 88000.0,
 8.7: 83000.0,
 10: 83000.0}
```

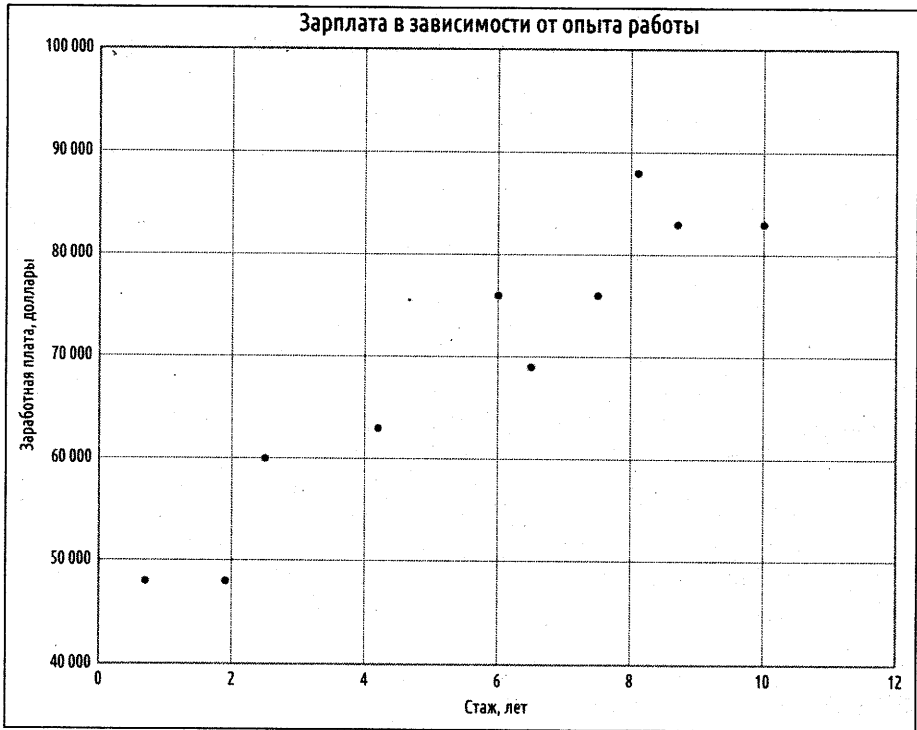


Рис. 1.3. Зависимость заработной платы от опыта работы

Целесообразнее разбить продолжительности стажа на интервалы:

```
# стажная группа
def tenure_bucket(tenure):
    if tenure < 2:
        return "менее двух"
    elif tenure < 5:
        return "между двумя и пятью"
    else:
        return "более пяти"
```

Затем сгруппировать все зарплаты, которые соответствуют каждому интервалу:

```
# зарплата в зависимости от стажной группы
# ключи = стажные группы, значения = списки зарплат в этой группе
# словарь содержит списки зарплат, соответствующие каждой стажной группе
salary_by_tenure_bucket = defaultdict(list)
```

```
for salary, tenure in salaries_and_tenures:
    bucket = tenure_bucket(tenure)
    salary_by_tenure_bucket[bucket].append(salary)
```

И в конце вычислить среднюю арифметическую зарплату по каждой группе:

```
# средняя зарплата по группе
# ключи = стажные группы, значения = средняя зарплата по этой группе
average_salary_by_bucket = {
    tenure_bucket : sum(salaries) / len(salaries)
    for tenure_bucket, salaries in salary_by_tenure_bucket.iteritems()
}
```

Результат выглядит гораздо интереснее:

```
{'между двумя и пятью': 61500.0,
 'менее двух': 48000.0,
 'более пяти': 79166.66666666667}
```

К тому же вы получаете короткий и запоминающийся "звуковой укус"<sup>4</sup>: "Аналитики с опытом работы свыше 5 лет зарабатывают на 65% больше, чем аналитики с малым или отсутствующим опытом работы в этой области!"

Впрочем, мы сегментировали данные на интервалы произвольным образом. На самом деле, хотелось бы получить некую констатацию степени воздействия на уровень зарплаты — в среднем — дополнительного года работы. Кроме получения более энергичного примечательного факта, это позволило бы строить предсказания о неизвестных зарплатах. Эта идея будет исследована в *главе 14*.

## Оплата премиум-аккаунтов

Когда вы возвращаетесь к своему рабочему столу, директор по доходам уже вас ожидает. Она желает лучше разобраться в том, какие пользователи оплачивают привилегированные аккаунты, а какие нет. (Она знает их имена, но эти сведения не имеют особой оперативной ценности.)

Вы обнаруживаете, что, по всей видимости, между опытом работы и оплатой аккаунтов есть некое соответствие:

---

<sup>4</sup> Звуковой укус (англ. *sound bite*) — эффектная короткая реплика. Согласно исследованиям, способность публики сосредоточиваться на том, что она слышит или видит, ограничена средним временем внимания, равным 30 с. Поэтому вся реклама говорит "звуковыми укусами", чтобы максимально воспользоваться ограниченным сроком внимания публики. — *Прим. пер.*

0.7 оплачено  
 1.9 не оплачено  
 2.5 оплачено  
 4.2 не оплачено  
 6 не оплачено  
 6.5 не оплачено  
 7.5 не оплачено  
 8.1 не оплачено  
 8.7 оплачено  
 10 оплачено

Пользователи с малой и большой продолжительностью опыта работы, как правило, аккаунты оплачивают, а пользователи со средней продолжительностью — нет.

Следовательно, если создать модель — хотя этих данных определенно недостаточно для того, чтобы строить на них модель — то можно попытаться спрогнозировать, что пользователи с малым и большим опытом работы оплатят аккаунты, а середняки — нет:

```
# предсказать платежи, исходя из стажа
def predict_paid_or_unpaid(years_experience):
    if years_experience < 3.0:
        return "оплачено"
    elif years_experience < 8.5:
        return "не оплачено"
    else:
        return "оплачено"
```

Разумеется, здесь проанализированы лишь точки отсечения.

Имея в распоряжении больше данных (и больше познаний в области математики), можно было бы построить модель, предсказывающую правдоподобие оплаты пользователем аккаунта, исходя из его опыта работы. Мы займемся исследованием этого типа задач в *главе 16*.

## Популярные темы

В тот момент, когда ваш первый рабочий день подходит к концу, директор по стратегии содержательного наполнения запрашивает у вас данные о том, какими темами пользователи интересуются больше всего, чтобы соответствующим образом распланировать свой календарь работы с электронным журналом. У вас уже имеются необработанные данные из проекта системы рекомендации новых друзей:

```
interests = [
    (0, "Hadoop"), (0, "Big Data"), (0, "HBase"), (0, "Java"),
    (0, "Spark"), (0, "Storm"), (0, "Cassandra"),
    (1, "NoSQL"), (1, "MongoDB"), (1, "Cassandra"), (1, "HBase"),
    (1, "Postgres"), (2, "Python"), (2, "scikit-learn"), (2, "scipy"),
    (2, "numpy"), (2, "statsmodels"), (2, "pandas"), (3, "R"), (3, "Python"),
    (3, "statistics"), (3, "regression"), (3, "probability"),
```

```
(4, "machine learning"), (4, "regression"), (4, "decision trees"),
(4, "libsvm"), (5, "Python"), (5, "R"), (5, "Java"), (5, "C++"),
(5, "Haskell"), (5, "programming languages"), (6, "statistics"),
(6, "probability"), (6, "mathematics"), (6, "theory"),
(7, "machine learning"), (7, "scikit-learn"), (7, "Mahout"),
(7, "neural networks"), (8, "neural networks"), (8, "deep learning"),
(8, "Big Data"), (8, "artificial intelligence"), (9, "Hadoop"),
(9, "Java"), (9, "MapReduce"), (9, "Big Data")
```

]

Простой, если не самый впечатляющий, способ найти наиболее популярные темы — просто подсчитать количества вхождений слов или их частотности<sup>5</sup>:

1. Перевести каждую тему в строчные буквы (пользователи могут напечатать их прописными).
2. Разбить их на слова.
3. Подсчитать результаты.

В коде это выглядит так:

*# слова и частотности*

```
words_and_counts = Counter(word
                             for user, interest in interests
                             for word in interest.lower().split())
```

Это позволяет легко перечислить слова, которые встречаются более одного раза:

```
for word, count in words_and_counts.most_common():
    if count > 1:
        print(word, count)
```

В итоге получаем ожидаемый результат (если только вы не планируете для разбиения строки на два слова использовать библиотеку "scikit-learn", которая в этом случае ожидаемых результатов не гарантирует):

```
learning 3
java 3
python 3
big 3
data 3
hbase 2
regression 2
cassandra 2
statistics 2
probability 2
hadoop 2
networks 2
```

---

<sup>5</sup> Термин "частотность" используется в более общем смысле, как показатель, выражающий собой число повторений или возникновения событий (процессов). — *Прим. пер.*

machine 2  
neural 2  
scikit-learn 2  
r 2

Мы обратимся к более продвинутым способам извлечения тематики из данных в *главе 20*.

## **Вперед**

Это был успешный день! Уставший, вы ускользаете из здания, прежде чем кому-нибудь еще удастся задать вам вопрос по поводу чего-нибудь еще. Хорошенько отоспитесь, потому что на завтра намечена программа ориентации для новых сотрудников. (Да-да, вы проработали целый день, даже не пройдя программы ориентации! Поставьте этот вопрос перед кадровой службой.)

# Интенсивный курс языка Python

Уже 25 лет, как народ сходит с ума от Пайтона, чему я не перестаю удивляться.

*Майкл Пейлин<sup>1</sup>*

Все новые сотрудники DataSciencester обязаны проходить программу ориентации персонала, самой интересной частью которой является интенсивный курс языка программирования Python.

Это не всеобъемлющее руководство по языку. Напротив, интенсивный курс предназначен для того, чтобы выделить те элементы языка, которые будут наиболее важны в дальнейшем (некоторым из них зачастую уделяют скромное внимание даже в учебниках).

## ОСНОВЫ

### Установка

Python можно скачать с [python.org](https://python.org). Однако если он еще не установлен, то вместо него рекомендуем дистрибутивный пакет Anaconda (<https://store.continuum.io/cshop/anaconda/>), который уже включает в себя большинство библиотек, необходимых для работы в области науки о данных.

На момент написания книги последняя версия языка была Python 3.4. В DataSciencester, тем не менее, используется старый надежный Python 2.7. Python 3 не совместим с Python 2, а многие важные библиотеки работают только с версией 2.7. Сообщество аналитиков данных по-прежнему прочно сидит на версии 2.7, и мы не будем исключением. Убедитесь, что получили именно эту версию<sup>2</sup>.

Если вы не используете дистрибутив Anaconda, то не забудьте установить менеджер пакетов pip (<https://pypi.python.org/pypi/pip>), позволяющий легко устанавли-

---

<sup>1</sup> Майкл Пейлин (1943) — британский актер, писатель, телеведущий, участник комик-группы, известной благодаря юмористическому телешоу "Летающий цирк Монти Пайтона", выходявшему на BBC в 1970-х, в честь которого был назван язык Python (см. [https://ru.wikipedia.org/wiki/Пейлин,\\_Майкл](https://ru.wikipedia.org/wiki/Пейлин,_Майкл)). — *Прим. пер.*

<sup>2</sup> Сегодня все основные библиотеки для анализа и визуализации данных прекрасно работают в Python 3, и при этом в нем отсутствуют проблемы с кириллическими шрифтами. С учетом этого в настоящей книге за основу взят Python 3, то есть все приведенные в книге примеры были протестированы в Python версии 3.5.2, форма записи приведена в соответствие с синтаксисом Python версии 3, и все редкие случаи расхождения в реализации задокументированы в сносках. — *Прим. пер.*

вать сторонние пакеты, поскольку некоторые из них нам понадобятся. Стоит также установить намного более удобную для работы интерактивную оболочку IPython (<http://ipython.org>).

(Дистрибутив Anaconda идет вместе с pip и IPython.)

Просто выполните:

```
pip install ipython
```

и бороздите Интернет в поисках решений по поводу того или иного загадочного сообщения об ошибке, вызываемого этой командой.

## Дзен языка Python

В основе языка лежат несколько *кратких принципов конструирования программ* (<http://legacy.python.org/dev/peps/pep-0020/>), именуемых "дзен языка Python". Их текст выдается интерпретатором, если ввести команду `import this`.

Один из самых обсуждаемых следующий:

"Должен существовать один — и, желательно, только один — очевидный способ сделать это".

Код, написанный в соответствии с этим "очевидным" способом (и, возможно, не совсем очевидным для новичка), часто характеризуется как "питоновский". Несмотря на то, что темой данной книги не является программирование на Python, мы будем иногда противопоставлять питоновские и непитоновские способы достигнуть одной и той же цели, но, как правило, предпочтение будет отдаваться питоновским способам.

## Пробельные символы

Во многих языках программирования для разграничения блоков кода используются фигурные скобки. В Python используются отступы:

```
# пример отступов во вложенных циклах for
for i in [1, 2, 3, 4, 5]:
    print(i)          # первая строка в блоке for i
    for j in [1, 2, 3, 4, 5]:
        print(j)      # первая строка в блоке for j
        print(i + j)  # последняя строка в блоке for j
    print(i)          # последняя строка в блоке for i
print("циклы закончились")
```

Это делает код легко читаемым, но в то же время заставляет следить за форматированием. Пробел внутри круглых и квадратных скобок игнорируется, что облегчает написание многословных выражений:

```
# пример многословного выражения
long_winded_computation = (1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9 + 10 +
                            11 + 12 + 13 + 14 + 15 + 16 + 17 + 18 + 19 + 20)
```



и легко читаемого кода:

```
# список списков
list_of_lists = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]

# такой список списков легче читается
easier_to_read_list_of_lists = [ [1, 2, 3],
                                  [4, 5, 6],
                                  [7, 8, 9] ]
```

Для продолжения оператора на следующей строке используется обратная косая черта, впрочем, такая запись будет применяться редко:

```
two_plus_three = 2 + \
                 3
```

В следствие форматирования кода пробельными символами возникают трудности при копировании и вставке кода в оболочку Python. Например, попытка скопировать следующий код:

```
for i in [1, 2, 3, 4, 5]:

    # обратите внимание на пустую строку
    print(i)
```

в стандартную оболочку Python вызовет ошибку:

```
# Ошибка нарушения отступа: ожидается блок с отступом
IndentationError: expected an indented block
```

потому что для интерпретатора пустая строка свидетельствует об окончании блока кода с циклом `for`.

Оболочка IPython располагает "волшебной" функцией `%paste`, которая правильно вставляет все то, что находится в буфере обмена, включая пробельные символы. Только одно это уже является веской причиной для того, чтобы использовать IPython.

## Модули

Некоторые библиотеки среды программирования на основе Python не загружаются по умолчанию. Это касается как функций, составляющих часть языка Python, так и сторонних инструментов, загружаемых самостоятельно. Для того чтобы эти инструменты можно было использовать, необходимо *импортировать* модули, которые их содержат.

Один из подходов заключается в том, чтобы просто импортировать сам модуль:

```
import re
my_regex = re.compile("[0-9]+", re.I)
```

Здесь `re` — это название модуля, содержащего функции и константы для работы с регулярными выражениями. Импортировав таким способом весь модуль, можно обращаться к функциям, предворяя их префиксом `re..`

Если в коде переменная с именем `re` уже есть, то можно воспользоваться псевдонимом модуля:

```
import re as regex
my_regex = regex.compile("[0-9]+", regex.I)
```

Псевдоним используют также в тех случаях, когда импортируемый модуль имеет громоздкое имя или когда в коде происходит частое обращение к модулю. Например, при визуализации данных на основе модуля `matplotlib` для него обычно используют следующий стандартный псевдоним:

```
import matplotlib.pyplot as plt
```

Если из модуля нужно получить несколько конкретных значений, то их можно импортировать в явном виде и использовать без ограничений:

```
from collections import defaultdict, Counter
lookup = defaultdict(int)
my_counter = Counter()
```

Естественно, можно, не мудрствуя лукаво, импортировать в пространство имен все содержимое модуля, но тогда возникает риск непреднамеренной перезаписи уже объявленных переменных:

```
match = 10
from re import * # в модуле re есть функция match
print(match)    # "<function re.match>"
```

Следует избегать такого рода простых решений.

## Арифметические операции

В Python версии 2.7 по умолчанию используется целочисленное деление<sup>3</sup>, благодаря чему выражение `5/2` равно `2`. В большинстве случаев требуется получить иной результат, поэтому файлы с кодом всегда следует начинать со следующей директивы:

```
from __future__ import division
```

Благодаря ей деление `5/2` даст `2.5`. Все примеры кода, приведенные в данной книге, содержат именно этот новый стиль деления. В ряде случаев, где требуется целочисленное деление, им можно воспользоваться при помощи двойной косой черты: `5//2`.

<sup>3</sup> Это замечание относится к Python 2, в котором целочисленное деление выполняется с округлением результата вниз до ближайшего целого, т. е. `5/2=2`. Чтобы изменить это поведение, можно поставить нули после точки (`5.0/2.0`), и в результате получим ожидаемый ответ `2.5`. В Python 3 целочисленное деление стало более понятным: `5/2` дает `2.5`, а чтобы выполнить округление, используют двойную косую: `5//2=2`. — *Прим. пер.*

## Функции

*Функция* — это правило, принимающее ноль или несколько входящих аргументов и возвращающее соответствующий результат. В Python функции обычно определяются при помощи оператора `def`:

```
def double(x):
    """здесь, когда нужно, размещают
    многострочный документирующий комментарий docstring,
    который поясняет, что именно функция вычисляет.
    Например, данная функция умножает входящее значение на 2"""
    return x * 2
```

Функции в Python рассматриваются как объекты *первого класса*. Это означает, что их можно присваивать переменным и передавать в другие функции так же, как любые другие аргументы:

```
# применить функцию f к единице
def apply_to_one(f):
    """вызывает функцию f с единицей в качестве аргумента"""
    return f(1)

my_double = double          # ссылка на ранее определенную функцию
x = apply_to_one(my_double) # = 2
```

Кроме того, можно легко создавать короткие *анонимные функции* или *лямбда-выражения*:

```
y = apply_to_one(lambda x: x + 4) # = 5
```

Лямбда-выражения можно присваивать переменным. Однако большинство программистов напротив рекомендуют пользоваться оператором `def`:

```
another_double = lambda x: 2 * x # так не делать
def another_double(x): return 2 * x # лучше так
```

Параметрам функции, помимо этого, можно передавать аргументы по умолчанию<sup>4</sup>, которые следует указывать только тогда, когда ожидается значение, отличающееся от значения по умолчанию:

```
def my_print(message="мое сообщение по умолчанию"):
    print(message)

my_print("привет") # напечатает 'привет'
my_print()        # напечатает 'мое сообщение по умолчанию'
```

Иногда целесообразно указывать аргументы по имени:

```
# функция вычитания
def subtract(a=0, b=0):
    return a - b
```

<sup>4</sup> Параметр указывается в объявлении функции и описывает значение, которое необходимо передать при ее вызове. Аргумент — это конкретное значение, передаваемое при вызове функции. — *Прим. пер.*

```

subtract(10, 5)    # возвращает 5
subtract(0, 5)    # возвращает -5
subtract(b=5)     # то же, что и в предыдущем примере

```

В дальнейшем функции будут использоваться очень часто.

## Строки

*Символьные строки* (или последовательности символов) с обеих сторон ограничиваются одинарными или двойными кавычками (они должны совпадать):

```

single_quoted_string = 'наука о данных' # одинарные
double_quoted_string = "наука о данных" # двойные

```

Обратная косая черта используется для кодирования специальных символов. Например:

```

tab_string = "\t"          # обозначает символ табуляции
len(tab_string)           # = 1

```

Если требуется непосредственно сама обратная косая черта, которая встречается в именах каталогов в операционной системе Windows или в регулярных выражениях, то при помощи `r""` можно создать *неформатированную строку*:

```

not_tab_string = r"\t"    # обозначает символы '\' и 't'
len(not_tab_string)      # = 2

```

Многострочные блоки текста создаются при помощи тройных одинарных (или двойных) кавычек:

```

multi_line_string = """Это первая строка.
это вторая строка
а это третья строка"""

```

## Исключения

Когда что-то идет не так, Python вызывает *исключение*. Необработанные исключения приводят к непредвиденному останову программы. Исключения обрабатываются при помощи операторов `try` и `except`:

```

try:
    print(0 / 0)
except ZeroDivisionError:
    print("нельзя делить на ноль")

```

Хотя во многих языках программирования использование исключений считается плохим стилем программирования, в Python нет ничего страшного, если они используются с целью сделать код чище, и мы будем иногда поступать именно так.

## Списки

Наверное, наиважнейшей структурой данных в Python является *список*. Это просто упорядоченная совокупность (или коллекция), похожая на массив в других языках программирования, но с дополнительными функциональными возможностями.

```
integer_list = [1, 2, 3] # список целых чисел
heterogeneous_list = ["строка", 0.1, True] # разнородный список
list_of_lists = [integer_list, heterogeneous_list, []] # список списков
```

```
list_length = len(integer_list) # длина списка = 3
list_sum = sum(integer_list) # сумма значений в списке = 6
```

Устанавливать значение и получать доступ к *n*-му элементу списка можно при помощи квадратных скобок<sup>5</sup>:

```
x = list(range(10)) # задает список [0, 1, ..., 9]
zero = x[0] # = 0, списки нуль-индексные, т. е. индекс 1-го элемента = 0
one = x[1] # = 1
nine = x[-1] # = 9, по-питоновски взять последний элемент
eight = x[-2] # = 8, по-питоновски взять предпоследний элемент
x[0] = -1 # теперь x = [-1, 1, 2, 3, ..., 9]
```

Помимо этого, квадратные скобки применяются для "нарезки" списков:

```
first_three = x[:3] # первые три = [-1, 1, 2]
three_to_end = x[3:] # с третьего до конца = [3, 4, ..., 9]
one_to_four = x[1:5] # с первого по четвертый = [1, 2, 3, 4]
last_three = x[-3:] # последние три = [7, 8, 9]
without_first_and_last = x[1:-1] # без первого и последнего = [1, 2, ..., 8]
copy_of_x = x[:] # копия списка x = [-1, 1, 2, ..., 9]
```

В Python имеется оператор `in`, который проверяет принадлежность элемента списку:

```
1 in [1, 2, 3] # True
0 in [1, 2, 3] # False
```

Проверка заключается в поочередном просмотре всех элементов, поэтому пользоваться им стоит только тогда, когда точно известно, что список небольшой или неважно, сколько времени уйдет на проверку.

Списки легко сцеплять друг с другом:

```
x = [1, 2, 3]
x.extend([4, 5, 6]) # теперь x = [1, 2, 3, 4, 5, 6]
```

<sup>5</sup> В отличие от Python 2, где `range` — это функция, которая возвращает список, в Python 3 `range` — это конструктор встроенного типа в виде иммутабельной последовательности, и, чтобы подключиться к функционалу обработки списков, его следует преобразовать в список, поэтому `list(range(10))` вместо `range(10)`. Еще одной главной особенностью Python 3 является то, что в операторе печати `print` следует использовать круглые скобки. — *Прим. пер.*

Если нужно оставить список `x` без изменений, то можно воспользоваться сложением списков:

```
x = [1, 2, 3]
y = x + [4, 5, 6]    # y = [1, 2, 3, 4, 5, 6]; x не изменился
```

Обычно к спискам добавляют по одному элементу за одну операцию:

```
x = [1, 2, 3]
x.append(0)        # теперь x = [1, 2, 3, 0]
y = x[-1]         # = 0
z = len(x)        # = 4
```

Нередко бывает удобно *распаковать* список, если известно, сколько элементов в нем содержится:

```
x, y = [1, 2]      # теперь x = 1, y = 2
```

Если с обеих сторон выражения число элементов не одинаково, то будет выдано сообщение об ошибке `ValueError`.

Для отбрасываемого значения обычно используется символ подчеркивания:

```
_, y = [1, 2]      # теперь y == 2, первый элемент не нужен
```

## Кортежи

*Кортежи* — это неизменяемые (или иммутабельные) двоюродные братья списков. Практически все, что можно делать со списком, не внося в него изменения, можно делать и с кортежем. Вместо квадратных скобок кортеж оформляют круглыми скобками, или вообще обходятся без них:

```
my_list = [1, 2]    # задать список
my_tuple = (1, 2)   # задать кортеж
other_tuple = 3, 4  # еще один кортеж
my_list[1] = 3      # теперь my_list = [1, 3]
```

**try:**

```
my_tuple[1] = 3
```

**except TypeError:**

```
print("кортеж изменять нельзя")
```

Кортежи обеспечивают удобный способ для возвращения из функций нескольких значений:

```
# функция возвращает сумму и произведение двух параметров
```

```
def sum_and_product(x, y):
```

```
    return (x + y), (x * y)
```

```
sp = sum_and_product(2, 3)    # = (5, 6)
```

```
s, p = sum_and_product(5, 10) # s = 15, p = 50
```

Кортежи (и списки) также используются во *множественном присваивании*:

```
x, y = 1, 2 # теперь x = 1, y = 2
x, y = y, x # обмен переменными по-питоновски; теперь x = 2, y = 1
```

## Словари

*Словарь* или ассоциативный список — это еще одна основная структура данных. В нем значения связаны с *ключами*, что позволяет быстро извлекать значение, соответствующее конкретному ключу:

```
empty_dict = {} # задать словарь по-питоновски
empty_dict2 = dict() # не совсем по-питоновски
grades = {"Joel" : 80, "Tim" : 95} # литерал словаря (оценки за экзамены)
```

Доступ к значению по ключу можно получить при помощи квадратных скобок:

```
joels_grade = grades["Joel"] # = 80
```

При попытке запросить значение, которое в словаре отсутствует, будет выдано сообщение об ошибке `KeyError`:

```
try:
    kates_grade = grades["Kate"]
except KeyError:
    print("оценки для Кэйт отсутствуют!")
```

Проверить наличие ключа можно при помощи оператора `in`:

```
joel_has_grade = "Joel" in grades # True
kate_has_grade = "Kate" in grades # False
```

Словари имеют метод `get()`, который при поиске отсутствующего ключа вместо вызова исключения возвращает значение по умолчанию:

```
joels_grade = grades.get("Joel", 0) # = 80
kates_grade = grades.get("Kate", 0) # = 0
no_ones_grade = grades.get("No One") # значение по умолчанию = None
```

Присваивание значения по ключу выполняется при помощи тех же квадратных скобок:

```
grades["Tim"] = 99 # заменяет старое значение
grades["Kate"] = 100 # добавляет третью запись
num_students = len(grades) # = 3
```

Словари часто используются в качестве простого способа представить структурные данные:

```
tweet = {
    "user" : "joelgrus",
    "text" : "Наука о данных - потрясающая тема",
    "retweet_count" : 100,
    "hashtags" : ["#data", "#science", "#datascience", "#awesome", "#yolo"]
}
```

Помимо поиска отдельных ключей можно обратиться ко всем сразу:

```
tweet_keys = tweet.keys() # список ключей
tweet_values = tweet.values() # список значений
tweet_items = tweet.items() # список кортежей (ключ, значение)

"user" in tweet_keys # True, но использует медленное in списка
"user" in tweet # по-питоновски, использует быстрое in словаря
"joelgrus" in tweet_values # True
```

Ключи должны быть неизменяемыми; в частности, в качестве ключей нельзя использовать списки. Если нужен составной ключ, то лучше воспользоваться кортежем или же найти способ, как преобразовать ключ в строку.

## Словарь `defaultdict`

Пусть в документе необходимо подсчитать слова. Очевидным решением задачи является создание словаря, в котором ключи — это слова, а значения — частотности слов (или количества вхождений слов в текст). Во время проверки слов в случае, если текущее слово уже есть в словаре, то его частотность увеличивается, а если отсутствует, то оно добавляется в словарь:

```
# частотности слов
word_counts = {}
document = {} # некий документ; здесь он пустой
for word in document:
    if word in word_counts:
        word_counts[word] += 1
    else:
        word_counts[word] = 1
```

Кроме этого, можно воспользоваться приемом под названием "лучше просить прощения, чем разрешения" и перехватывать ошибку при попытке обратиться к отсутствующему ключу:

```
word_counts = {}
for word in document:
    try:
        word_counts[word] += 1
    except KeyError:
        word_counts[word] = 1
```

Третий прием — использовать метод `get()`, который изящно выходит из ситуации с отсутствующими ключами:

```
word_counts = {}
for word in document:
    previous_count = word_counts.get(word, 0)
    word_counts[word] = previous_count + 1
```

Все перечисленные приемы немного громоздки, и по этой причине целесообразно использовать словарь `defaultdict` (который еще называют словарем со значением по



умолчанию). Он похож на обычный словарь за исключением одной особенности — при попытке обратиться к ключу, которого в нем нет, он сперва добавляет для него значение, используя функцию без аргументов, которая предоставляется при его создании. Чтобы воспользоваться словарями `defaultdict`, их необходимо импортировать из модуля `collections`:

```
from collections import defaultdict

word_counts = defaultdict(int)    # int() возвращает 0
for word in document:
    word_counts[word] += 1
```

Кроме того, использование словарей `defaultdict` имеет практическую пользу во время работы со списками, словарями и даже с пользовательскими функциями:

```
dd_list = defaultdict(list)    # list() возвращает пустой список
dd_list[2].append(1)           # теперь dd_list содержит {2: [1]}

dd_dict = defaultdict(dict)    # dict() возвращает пустой словарь dict
dd_dict["Joel"]["City"] = "Seattle" # { "Joel" : { "City" : Seattle" }}

dd_pair = defaultdict(lambda: [0, 0])
dd_pair[2][1] = 1              # теперь dd_pair содержит {2: [0,1]}
```

Эти возможности понадобятся, когда словари будут использоваться для "сбора" результатов по некоторому ключу и когда необходимо избежать повторяющихся проверок на присутствие ключа в словаре.

## Словарь `Counter`

Подкласс словарей `Counter` трансформирует последовательность значений в похожий на словарь `defaultdict(int)` объект, где ключам поставлены в соответствие частотности или, выражаясь более точно, ключи отображаются (`map`) в частотности. Он в основном будет применяться при создании гистограмм:

```
from collections import Counter
c = Counter([0, 1, 2, 0])    # в результате c = { 0 : 2, 1 : 1, 2 : 1 }
```

Его функционал позволяет достаточно легко решить задачу подсчета частотностей слов:

```
# лучший вариант подсчета частотностей слов
word_counts = Counter(document)
```

Словарь `Counter` располагает методом `most_common()`, который нередко бывает полезен:

```
# напечатать 10 наиболее встречаемых слов и их частотность (встречаемость)
for word, count in word_counts.most_common(10):
    print(word, count)
```

## Множества

Структура данных `set` или *множество* представляет собой совокупность неупорядоченных элементов *без повторов*:

```
s = set()      # задать пустое множество
s.add(1)      # теперь s = { 1 }
s.add(2)      # теперь s = { 1, 2 }
s.add(2)      # s как и прежде = { 1, 2 }
x = len(s)    # = 2
y = 2 in s    # = True
z = 3 in s    # = False
```

Множества будут использоваться по двум причинам. Во-первых, операция `in` на множествах очень быстрая. Если необходимо проверить большую совокупность элементов на принадлежность некоторой последовательности, то структура данных `set` подходит для этого лучше, чем список:

```
# список стоп-слов
stopwords_list = ["a", "an", "at"] + hundreds_of_other_words + ["yet", "you"]
"zip" in stopwords_list # False, но проверяется каждый элемент
```

```
# множество стоп-слов
stopwords_set = set(stopwords_list)
"zip" in stopwords_set # очень быстрая проверка
```

Вторая причина — получение *уникальных* элементов в наборе данных:

```
item_list = [1, 2, 3, 1, 2, 3] # список
num_items = len(item_list)    # количество = 6
item_set = set(item_list)     # вернет множество {1, 2, 3}
num_distinct_items = len(item_set) # число недублирующихся = 3
distinct_item_list = list(item_set) # назад в список = [1, 2, 3]
```

Множества будут применяться намного реже словарей и списков.

## Управляющие конструкции

Как и в большинстве других языков программирования, действия можно выполнять по условию, применяя оператор `if`:

```
if 1 > 2:
    message = "если 1 была бы больше 2..."
elif 1 > 3:
    message = "elif означает 'else if'"
else:
    message = "когда все предыдущие условия не выполняются, используется else"
```

Кроме того, можно воспользоваться однострочным *трехместным* оператором `if-then-else`, который будет иногда использоваться в дальнейшем:

```
parity = "четное" if x % 2 == 0 else "нечетное"
```

В Python имеется цикл `while`:

```
x = 0
while x < 10:
    print(x, "меньше 10")
    x += 1
```

Однако чаще будет использоваться цикл `for` совместно с оператором `in`:

```
for x in range(10):
    print(x, "меньше 10")
```

Если требуется более сложная логика управления циклом, то можно воспользоваться операторами `continue` и `break`:

```
for x in range(10):
    if x == 3:
        continue # перейти сразу к следующей итерации
    if x == 5:
        break # выйти из цикла
    print(x)
```

В результате будет напечатано 0, 1, 2 и 4.

## ИСТИННОСТЬ

Булевы переменные в Python работают так же, как и в большинстве других языков программирования лишь с одним исключением — они пишутся с заглавной буквы:

```
one_is_less_than_two = 1 < 2 # = True
true_equals_false = True == False # = False
```

Для обозначения несуществующего значения применяется специальный объект `None`, который соответствует значению `null` в других языках:

```
x = None
print(x == None) # напечатает True, но это не по-питоновски
print(x is None) # напечатает True по-питоновски
```

В Python может использоваться любое значение там, где ожидается логический тип `Boolean`. Все следующие элементы имеют логическое значение `False`:

- ◆ `False`;
- ◆ `None`;
- ◆ `[]` (пустой список);
- ◆ `{}` (пустой словарь);
- ◆ `""`;
- ◆ `set()` (множество);
- ◆ `0`;
- ◆ `0.0`.

Практически все остальное рассматривается как `True`. Это позволяет легко применять операторы `if` для проверок на наличие пустых списков, пустых строк, пустых словарей и т. д. Иногда, правда, это приводит к труднораспознаваемым ошибкам, если не учитывать следующее:

```
s = some_function_that_returns_a_string() # возвращает строковое значение
if s:
    first_char = s[0] # первый символ в строке
else:
    first_char = ""
```

Вот более простой способ сделать то же самое:

```
first_char = s and s[0]
```

поскольку логический оператор `and` возвращает второе значение, в случае если первое истинное, и первое значение, в случае если оно ложное. Аналогичным образом, если `x` в следующем ниже выражении является либо числом, либо, возможно, `None`, то результат так или иначе будет числом:

```
safe_x = x or 0 # безопасный способ
```

Встроенная функция `all` языка Python берет список и возвращает `True` только тогда, когда каждый элемент списка истинен, а встроенная функция `any` возвращает `True`, когда истинен хотя бы один элемент:

```
all([True, 1, { 3 }]) # True
all([True, 1, {}]) # False, {} = ложное
any([True, 1, {}]) # True, True = истинное
all([]) # True, ложные элементы в списке отсутствуют
any([]) # False, истинные элементы в списке отсутствуют
```

## Не совсем основы

В этом разделе мы обратимся к некоторым более сложным элементам функционала языка Python, которые будут полезны во время работы с данными.

### Сортировка

Каждый список в Python располагает методом `sort()`, который упорядочивает его на месте, т. е. внутри списка без выделения дополнительной памяти. Чтобы избежать изменений в списке, можно применить встроенную функцию `sorted`, которая возвращает новый список:

```
x = [4, 1, 2, 3]
y = sorted(x) # после сортировки = [1, 2, 3, 4], x не изменился
x.sort() # теперь x = [1, 2, 3, 4]
```

По умолчанию метод `sort()` и функция `sorted` сортируют список по возрастанию значению элементов, сравнивая элементы друг с другом.

Если необходимо, чтобы элементы были отсортированы по убывающему значению, надо указать аргумент направления сортировки `reverse=True`. А вместо сравнения самих элементов можно сравнивать результаты функции, которую надо указать вместе с ключом `key`:

```
# сортировать список по абсолютному значению в убывающем порядке
x = sorted([-4,1,-2,3], key=abs, reverse=True) # = [-4,3,-2,1]

# сортировать слова и их частотности по убывающему значению частот
wc = sorted(word_counts.items(),
            key=lambda word: count, reverse=True)
```

## Генераторы последовательностей

Нередко появляется необходимость преобразовать некий список в другой, выбирая только определенные элементы или внося по ходу изменения в элементы исходного списка, или выполняя и то и другое одновременно. Питоновским решением является применение *генераторов последовательностей* (еще именуемых списковыми включениями)<sup>6</sup>:

```
# четные числа
even_numbers = [x for x in range(5) if x % 2 == 0] # [0, 2, 4]

# квадраты чисел
squares = [x * x for x in range(5)] # [0, 1, 4, 9, 16]

# квадраты четных чисел
even_squares = [x * x for x in even_numbers] # [0, 4, 16]
```

Списки можно преобразовать в словари dict или множества set аналогичным образом:

```
# словарь с квадратами чисел
square_dict = { x : x * x for x in range(5) } # { 0:0, 1:1, 2:4, 3:9, 4:16 }

# множество с квадратами чисел
square_set = { x * x for x in [1, -1] } # { 1 }
```

Если значение в списке не нужно, то обычно в качестве переменной используется символ подчеркивания:

```
# нули
zeroes = [0 for _ in even_numbers] # имеет ту же длину, что и even_numbers
```

Генератор последовательности может содержать несколько операторов for:

```
# пары
pairs = [(x, y)
         for x in range(10)
         for y in range(10)] # 100 пар (0,0) (0,1)... (9,8), (9,9)
```

И последующие операторы for могут использовать результаты предыдущих:

```
# пары с возрастающим значением
increasing_pairs = [(x, y) # только пары с x < y,
                   for x in range(10) # range(мин, макс) равен
                   for y in range(x + 1, 10)] # [мин, мин+1, ..., макс-1]
```

Генераторы последовательностей будут применяться очень часто.

<sup>6</sup> Генератор последовательности аналогичен математической записи для задания множества путем описания свойств, которыми должны обладать его члены — *Прим. пер.*

## Функции-генераторы и генераторные выражения

Проблема со списками заключается в том, что они легко вырастают до больших размеров. Выражение `range(1000000)` создает фактическую последовательность (в Python 2 это выражение создает список) из 1 млн элементов. Если их обрабатывать по одному за итерацию цикла, то они могут стать главной причиной неэффективности (или нехватки оперативной памяти). А если потенциально будут использоваться лишь первые несколько значений, то вычисление всех элементов становится пустой тратой ресурсов.

*Генератор* — это объект, который можно последовательно перебрать (обычно при помощи оператора `for`), но чьи значения предоставляются только тогда, когда они требуются (используя *ленивое вычисление*).

Первый способ создания генераторов заключается в использовании функции-генератора совместно с оператором `yield`:

```
def lazy_range(n):
    """ленивая версия последовательности range"""
    i = 0
    while i < n:
        yield i
        i += 1
```

Представленный далее цикл будет потреблять предоставляемые значения по одному за итерацию, пока элементы не закончатся:

```
for i in lazy_range(10): do_something_with(i)
```

(На самом деле Python уже располагает функцией `lazy_range` под названием `xrange`, а в версии Python 3 сама последовательность (диапазон `range`) уже "ленивая".) Такая функциональность позволяет задавать бесконечную последовательность:

```
# натуральные числа
def natural_numbers():
    """возвращает 1, 2, 3, ..."""
    n = 1
    while True:
        yield n
        n += 1
```

Хотя, наверное, не стоит просматривать такую последовательность без применения какой-нибудь схемы выхода из цикла.



Обратной стороной ленивого вычисления является то, что просмотр генератора можно выполнить всего один раз. Если нужно просмотреть многократно, то следует либо каждый раз создавать генератор заново, либо использовать список.

Второй способ создания генераторов — использовать генераторное выражение, т. е. генератор последовательности с оператором `for`, обрамленный круглыми скобками:

```
# ленивый список четных чисел меньше 20
lazy_evens_below_20 = (i for i in lazy_range(20) if i % 2 == 0).
```

Следует отметить, что каждый объект dict имеет метод items(), который возвращает список пар "ключ-значение". Однако чаще будет применяться метод iteritems(), который "лениво" предоставляет по одной паре по мере навигации по словарю:

## Случайные числа

При изучении науки о данных часто возникает необходимость генерировать случайные числа. Для этого предназначен модуль random:

```
import random

# четыре равномерные случайные величины
four_uniform_randoms = [random.random() for _ in range(4)]

# [0.8444218515250481,      # random.random() производит числа
# 0.7579544029403025,      # равномерно в интервале между 0 и 1
# 0.420571580830845,      # функция random будет применяться
# 0.25891675029296335]    # наиболее часто
```

Модуль random на самом деле генерирует псевдослучайные (т. е. детерминированные) числа на основе внутреннего состояния. Когда требуется получить воспроизводимые результаты, внутреннее состояние можно задать при помощи метода random.seed:

```
random.seed(10)      # задать случайную последовательность, установив в 10
print(random.random()) # 0.57140259469
random.seed(10)      # переустановить seed в 10
print(random.random()) # опять 0.57140259469
```

Иногда будет применяться метод random.randrange, который принимает один или два аргумента и возвращает элемент, выбранный случайно из соответствующей последовательности range():

```
random.randrange(10) # произвольно выбрать из range(10) = [0, 1, ..., 9]
random.randrange(3, 6) # произвольно выбрать из range(3, 6) = [3, 4, 5]
```

Есть еще ряд методов, которые пригодятся. Метод random.shuffle, например, перемешивает элементы в списке в случайном порядке:

```
up_to_ten = range(10) # задать последовательность из 10 элементов
random.shuffle(up_to_ten)
print(up_to_ten)
# [2, 5, 1, 9, 7, 3, 8, 6, 4, 0] (фактические результаты могут отличаться)
```

Если нужно случайным образом выбрать из списка один элемент, то можно воспользоваться методом random.choice:

```
# мой лучший друг
my_best_friend = random.choice(["Алиса", "Боб", "Чарли"]) # сейчас "Боб"
```

Если из исходных данных необходимо сделать произвольную выборку элементов *без возврата* (т. е. без повторяющихся элементов; этот способ отбора еще называется *безповторным*), то для этих целей служит метод `random.sample`:

```
# лотерейные номера
lottery_numbers = range(60)
# выигрышные номера (пример выборки без возврата)
winning_numbers = random.sample(lottery_numbers, 6) # [16, 36, 10, 6, 25, 9]
```

Для получения выборки элементов *с возвратом* (т. е. допускающую повторы; этот способ отбора еще называется *повторным*) надо всего лишь сделать несколько вызовов метода `random.choice()`:

```
# список из четырех элементов (пример выборки с возвратом)
four_with_replacement = [random.choice(range(10))
                          for _ in range(4)]
# [9, 4, 4, 2]
```

## Регулярные выражения

Регулярные выражения предоставляют средства для выполнения поиска в тексте. Они невероятно полезны и одновременно довольно сложны, причем настолько, что им посвящены целые книги. В дальнейшем в тех нескольких случаях, когда они будут встречаться в коде, будут даваться разъяснения, связанные с их работой. Вот несколько примеров их применения в Python:

```
import re

print(all([
    not re.match("a", "cat"),          # все они - истинные, т. к.
    re.search("a", "cat"),             # слово 'cat' не начинается с 'a'
    not re.search("c", "dog"),        # в слове 'cat' есть 'a'
    3 == len(re.split("[ab]", "carbs")), # в слове 'dog' нет 'c'
                                        # разбивка по a или b
                                        # в ['c', 'r', 's']
    "R-D-" == re.sub("[0-9]", "-", "R2D2") # замена цифр дефисами
])) # напечатает True
```

## Объектно-ориентированное программирование

Как и во многих языках программирования, в Python имеется возможность определять *классы*, которые содержат данные и функции для их обработки. Мы иногда будем ими пользоваться для того, чтобы сделать код чище и проще. Наверное, самый простой способ объяснить, как они работают — привести пример с подробным описанием.

Представим, что в Python отсутствует встроенная структура данных `Set` для реализации множеств. Тогда может понадобиться создать собственный класс.

Каким поведением этот класс должен обладать? *Экземпляр* класса `Set` должен уметь добавлять (`add`) и удалять (`remove`) элементы, а также проверять, содержит ли он



определенное значение (`contains`). Эту функциональность можно задать в виде *компонентных функций*, доступ к которым обеспечивается через точку после идентификатора объекта `Set`:

```
# по традиции классам назначают имена с заглавной буквы
class Set:

    # ниже идут компонентные функции
    # каждая берет первый параметр "self" (еще одно правило),
    # который ссылается на конкретный используемый объект класса Set
    def __init__(self, values=None):
        """Это конструктор.
        Вызывается при создании нового объекта класса Set и
        используется следующим образом:
        s1 = Set()           # пустое множество
        s2 = Set([1,2,2,3]) # инициализировать значениями"""

        self.dict = {} # каждый экземпляр имеет собственное свойство dict,
                       # которое используется для проверки
                       # на принадлежность элементов множеству

        if values is not None:
            for value in values: self.add(value)

    def __repr__(self):
        """это строковое представление объекта Set,
        которое выводится в оболочке или передается в функцию str()"""
        return "Set: " + str(self.dict.keys())

    # принадлежность представлена ключом в словаре self.dict
    # со значением True
    def add(self, value):
        self.dict[value] = True

    # значение принадлежит множеству, если его ключ имеется в словаре
    def contains(self, value):
        return value in self.dict

    def remove(self, value):
        del self.dict[value]
```

Все это можно использовать следующим образом:

```
s = Set([1,2,3])
s.add(4)
print(s.contains(4))    # True
s.remove(3)
print(s.contains(3))    # False
```

## Инструменты функционального программирования

Передавая функции в качестве аргументов, иногда требуется, чтобы они применялись частично (или *каррировались*) с целью создания новых функций. В качестве простого примера, пусть имеется функция двух переменных:

```
# возведение в степень power
def exp(base, power):
    return base ** power
```

И надо применить ее для создания функции одной переменной `two_to_the` для возведения в степень 2, чьим входящим аргументом является функция `power`, а выходящим значением — результат вычисления функции `exp(2, power)`.

Разумеется, эту задачу можно решить при помощи оператора `def`, но иногда это становится неудобным:

```
# двойка в степени power
def two_to_the(power):
    return exp(2, power)
```

Другой прием заключается в использовании метода `functools.partial`:

```
from functools import partial
# возвращает функцию с частичным приложением аргументов
two_to_the = partial(exp, 2) # теперь это функция одной переменной
print(two_to_the(3))       # 8
```

Метод `partial()` также используют для заполнения последующих аргументов, указывая при этом их имена:

```
square_of = partial(exp, power=2) # квадрат числа
print(square_of(3))              # 9
```

Часто возникает путаница, если каррировать аргументы в середине функции. Следует избегать такого приема.

В дальнейшем будут периодически использоваться встроенные функции `map`, `reduce` и `filter`, которые предоставляют функциональные альтернативы генераторам последовательностей:

```
def double(x):
    return 2 * x

xs = [1, 2, 3, 4]
twice_xs = [double(x) for x in xs] # [2, 4, 6, 8]
twice_xs = map(double, xs)       # то же, что и выше
list_doubler = partial(map, double) # удвоитель списка
twice_xs = list_doubler(xs)      # снова [2, 4, 6, 8]
```

Функцию `map` используют с функциями нескольких аргументов, предоставляя ей несколько списков:

```
# перемножить аргументы
def multiply(x, y): return x * y

products = map(multiply, [1, 2], [4, 5]) # [1 * 4, 2 * 5] = [4, 10]
```

Функция `filter` работает аналогично генератору последовательности с условием `if`:

```
# проверка четности
def is_even(x):
    """True, если x - четное; False, если x - нечетное"""
    return x % 2 == 0

x_evens = [x for x in xs if is_even(x)] # список четных чисел = [2, 4]
x_evens = filter(is_even, xs) # то же, что и выше
list_evener = partial(filter, is_even) # функция, которая фильтрует
# список
x_evens = list_evener(xs) # снова [2, 4]
```

Функция `reduce` выполняет свертку списка, т. е. объединяет первые два элемента, затем этот результат с третьим элементом, а тот с четвертым и т. д., возвращая единственный результат:

```
from functools import reduce

x_product = reduce(multiply, xs) # = 1 * 2 * 3 * 4 = 24
list_product = partial(reduce, multiply) # функция, которая упрощает
# список
x_product = list_product(xs) # снова = 24
```

## Функция `enumerate`

Нередко нужно итеративно обойти список, используя как элементы, так и их индексы:

```
documents = [] # список неких документов; здесь он пустой
# не по-питоновски
for i in range(len(documents)):
    document = documents[i]
    do_something(i, document)

# тоже не по-питоновски
i = 0
for document in documents:
    do_something(i, document)
    i += 1
```

Питоновским решением является встроенная функция `enumerate`, которая создает кортежи в формате (индекс, элемент):

```
for i, document in enumerate(documents):
    do_something(i, document)
```

Точно так же, если нужны только индексы:

```
for i in range(len(documents)): do_something(i)      # не по-ПИТОНОВСКИ
for i, _ in enumerate(documents): do_something(i)   # по-ПИТОНОВСКИ
```

Все это будет широко использоваться в дальнейшем.

## Функция *zip* и распаковка аргументов

Часто возникает необходимость объединить два или более списков в один. Встроенная функция *zip* преобразует несколько списков в один список кортежей, состоящий из соответствующих элементов:

```
list1 = ['a', 'b', 'c']
list2 = [1, 2, 3]
list(zip(list1, list2))    # = [('a', 1), ('b', 2), ('c', 3)]
```

Если списки имеют разные длины, то функция прекратит работу, как только закончится первый из них. "Разъединить" список можно при помощи замысловатого трюка:

```
pairs = [('a', 1), ('b', 2), ('c', 3)]
letters, numbers = zip(*pairs)
```

Звездочка *\** выполняет *распаковку аргументов*, которая использует элементы пар в качестве индивидуальных аргументов для функции *zip*. Это то же самое, если вызвать:

```
list(zip(('a', 1), ('b', 2), ('c', 3)))
```

Результат в обоих случаях будет [('a', 'b', 'c'), (1, 2, 3)].

Распаковку аргументов можно применять с любой функцией:

```
def add(a, b): return a + b

add(1, 2)      # вернет 3
add([1, 2])   # ошибка TypeError!
add(*[1, 2])  # вернет 3
```

Этот прием будет использоваться редко; и то в качестве программистского трюка.

## Переменные *args* и *kwargs*

Предположим, нужно создать функцию высшего порядка, которая в качестве входящего аргумента принимает некую функцию *f* и возвращает новую функцию, которая для любого входящего аргумента возвращает удвоенное значение *f*:

```
# удвоитель
def doubler(f):
    def g(x):
        return 2 * f(x)
    return g
```

В некоторых случаях это работает:

```
def f1(x):
    return x + 1

g = doubler(f1)
print(g(3))      # 8 или (3 + 1) * 2
print(g(-1))    # 0 или (-1 + 1) * 2
```

Но терпит неудачу с функциями, которые принимают более одного аргумента:

```
def f2(x, y):
    return x + y

g = doubler(f2)
print(g(1, 2))  # TypeError: g() принимает ровно 1 аргумент (задано 2)
```

Хотелось бы определить функцию, которая принимает произвольное количество аргументов. Это делается при помощи распаковки аргументов и небольшого волшебства:

```
def magic(*args, **kwargs):
    print("безымянные аргументы:", args)
    print("аргументы по ключу:", kwargs)

magic(1, 2, key="word", key2="word2")

# напечатает
# безымянные аргументы: (1, 2)
# аргументы по ключу: {'key2': 'word2', 'key': 'word'}
```

Другими словами, когда функция определена подобным образом, переменная `args` представляет собой кортеж из безымянных позиционных аргументов, а переменная `kwargs` — словарь из именованных аргументов. Она работает и по-другому, если нужно воспользоваться списком (или кортежем) и словарем для того, чтобы *передать* аргументы в функцию:

```
def other_way_magic(x, y, z):
    return x + y + z

x_y_list = [1, 2]
z_dict = {"z": 3}
print(other_way_magic(*x_y_list, **z_dict))  # 6
```

Эти свойства переменных `args` и `kwargs` позволяют проделывать разного рода замысловатые трюки. Здесь этот прием будет использоваться исключительно для создания функций высшего порядка, которые могут принимать произвольное число входящих аргументов:

```
# корректный удвоитель
def doubler_correct(f):
```

```
"""работает независимо от того, какого рода аргументы
   функция f ожидает"""
def g(*args, **kwargs):
    """какими бы ни были аргументы для g, передать их в f"""
    return 2 * f(*args, **kwargs)
return g

g = doubler_correct(f2)
print(g(1, 2)) # 6
```

## Добро пожаловать в DataSciencester!

На этом программа ориентации нового персонала завершается. Да, кстати, постарайтесь не растерять полученные сведения.

## Для дальнейшего изучения

- ◆ Дефицит пособий по языку Python в мире отсутствует. Официальное руководство (<https://docs.python.org/2/tutorial/>) будет неплохим пособием для начала.
- ◆ Официальное руководство по IPython (<http://ipython.org/ipython-doc/2/interactive/tutorial.html>) не так удачно. Лучше ознакомиться с видеоуроками и презентациями (<http://ipython.org/videos.html>). Как вариант, в книге Уэса Маккинни "Python и анализ данных"<sup>7</sup> (<http://shop.oreilly.com/product/0636920023784.do>) имеется действительно хорошая статья по IPython.

---

<sup>7</sup> Маккинни У. Python и анализ данных. — М.: ДМК Пресс, 2015.

# Визуализация данных

Думаю, что визуализация — одно из самых мощных орудий достижения личных целей.

*Харви Маккей<sup>1</sup>*

Важнейшую часть инструментария аналитика данных составляют графические методы анализа данных или методы визуализации. И если создать визуализацию достаточно просто, то сделать *хорошую* визуализацию — задача не из легких.

Визуализация применяется прежде всего для:

- ◆ *исследования* данных;
- ◆ *передачи* данных.

Основная задача этой главы — развитие навыков, необходимых для того, чтобы вы могли начать исследование собственных данных и создавать на их основе визуализации, которые мы будем использовать на протяжении всей книги. Как и большинство тем, вынесенных в названия глав этой книги, визуализация данных представляет собой обширную область для изучения, которая заслуживает отдельной книги. Тем не менее, мы постараемся дать вам общее представление о том, что такое хорошая визуализация.

## Библиотека `matplotlib`

Для графического анализа данных существует большой набор разных инструментов. Мы будем пользоваться популярной библиотекой `matplotlib` (<http://matplotlib.org/>) (которая, впрочем, уже зарекомендовала себя в качестве хорошего инструмента). Если вас интересует создание детально проработанных интерактивных визуализаций для Интернета, то эта библиотека, скорее всего, вам не подойдет. Но с построением элементарных столбчатых и точечных диаграмм и линейных графиков она справляется легко.

В частности, мы будем пользоваться модулем `matplotlib.pyplot`. В самом простом случае `pyplot` позволяет выстраивать визуализацию шаг за шагом. Когда она завершена, изображение можно сохранить в файл (при помощи метода `savefig()`) или вывести на экран (методом `show()`).

---

<sup>1</sup> Харви Маккей (1932) — предприниматель, писатель и обозреватель крупнейших американских новостных изданий (см. [https://en.wikipedia.org/wiki/Harvey\\_Mackay](https://en.wikipedia.org/wiki/Harvey_Mackay)). — *Прим. пер.*

Например, простой график, такой как на рис. 3.1, создается достаточно легко:

```
from matplotlib import pyplot as plt
```

```
years = [1950, 1960, 1970, 1980, 1990, 2000, 2010] # годы
```

```
gdp = [300.2, 543.3, 1075.9, 2862.5, 5979.6, 10289.7, 14958.3] # ВВП
```

```
# создать линейную диаграмму: годы по оси X, ВВП по оси Y
```

```
plt.plot(years, gdp, color='green', marker='o', linestyle='solid')
```

```
# добавить название диаграммы
```

```
plt.title("Номинальный ВВП")
```

```
# добавить подпись к оси Y
```

```
plt.ylabel("Млрд $")
```

```
plt.show()
```

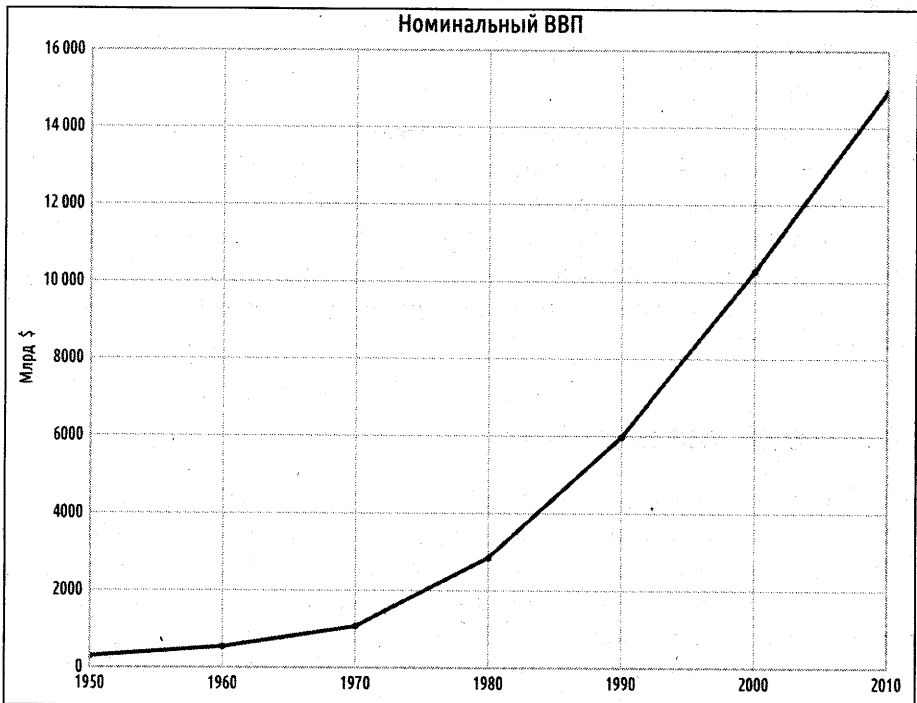


Рис. 3.1. Простая линейная диаграмма

Создание диаграмм типографского качества представляет собой более сложную задачу и выходит за рамки этой главы. Вы можете настраивать диаграммы самыми разнообразными способами, например, задавая подписи к осям, типы линий, маркеры точек и т. д. Вместо того, чтобы подробно изучать все эти возможности, мы просто будем пользоваться некоторыми из них в наших примерах и разбирать особенности их реализации по ходу изложения.





Значительная часть функционала библиотеки `matplotlib` в книге не используется. Тем не менее, надо отметить, что библиотека позволяет создавать составные вложенные диаграммы, выполнять сложное форматирование и делать интерактивные визуализации. Сверьтесь с документацией, если хотите копнуть глубже, чем это представлено в книге.

## Столбчатые диаграммы

Столбчатые диаграммы хорошо подходят для тех случаев, когда требуется показать *варьируемость* некоторой величины на некотором дискретном множестве элементов. Например, на рис. 3.2 показано, сколько ежегодных премий "Оскар" Американской киноакадемии было завоевано каждым из перечисленных фильмов:

```
movies = ["Энни Холл", "Бен-Гур", "Касабланка", "Ганди", "Вестсайдская история"]
num_oscars = [5, 11, 3, 8, 10]
```

```
# ширина столбцов по умолчанию 0.8, поэтому добавим 0.1 к левым
# координатам, чтобы каждый столбец был по центру интервала
xs = [i + 0.1 for i, _ in enumerate(movies)]

# построить столбцы с левыми X-координатами [xs] и высотой [num_oscars]
plt.bar(xs, num_oscars)
```

```
plt.ylabel("Количество наград")
plt.title("Мои любимые фильмы")
```

```
# добавить метки на оси X с названиями фильмов в центре каждого интервала
plt.xticks([i + 0.5 for i, _ in enumerate(movies)], movies)
plt.show()
```

Столбчатая диаграмма также хорошо подходит для построения гистограмм сгруппированных значений числового ряда, что позволяет наглядно исследовать характер *распределения* чисел в числовом ряду (рис. 3.3):

```
grades = [83, 95, 91, 87, 70, 0, 85, 82, 100, 67, 73, 77, 0]
decile = lambda grade: grade // 10 * 10 # Дециль (десятая часть числа)
histogram = Counter(decile(grade) for grade in grades)
```

```
plt.bar([x - 4 for x in histogram.keys()], # сдвинуть столбец влево на 4
        histogram.values(),             # высота столбца
        8)                               # ширина каждого столбца 8
```

```
plt.axis([-5, 105, 0, 5])               # ось X от -5 до 105,
                                         # ось Y от 0 до 5
```

```
plt.xticks([10 * i for i in range(11)]) # метки по оси X: 0, 10, ..., 100
plt.xlabel("Дециль")
plt.ylabel("Число студентов")
plt.title("Распределение оценок за экзамен № 1")
plt.show()
```

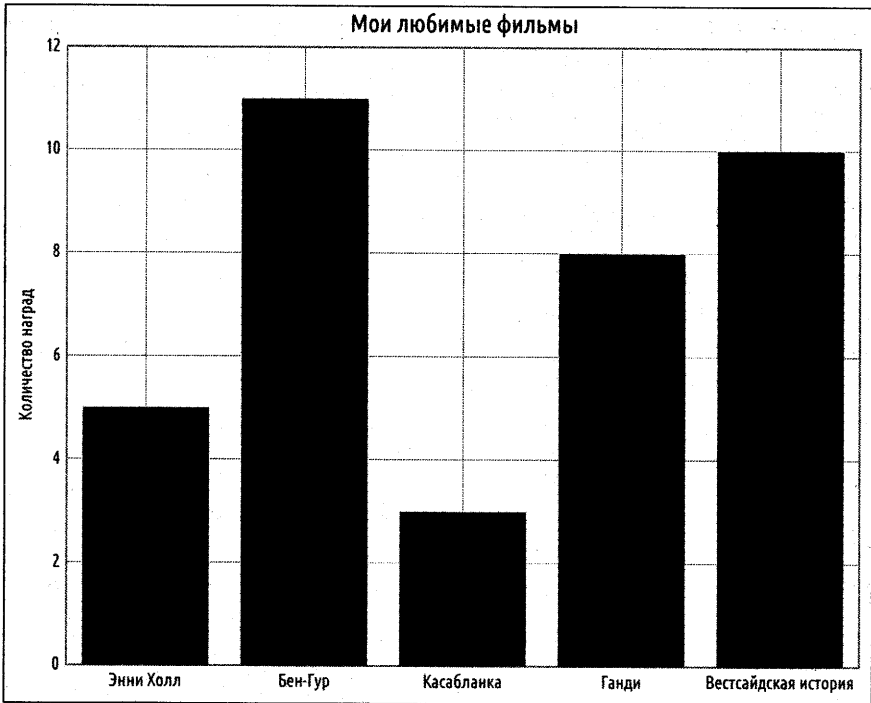


Рис. 3.2. Простая столбчатая диаграмма

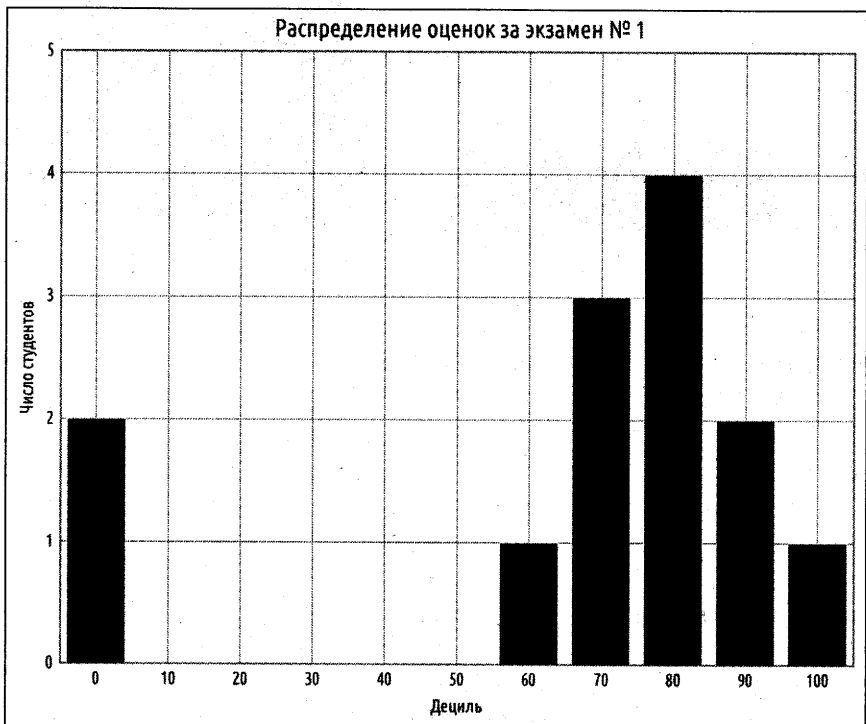


Рис. 3.3. Использование столбчатой диаграммы для гистограммы

Третий аргумент в методе `plt.bar` определяет ширину столбцов. На нашей диаграмме ширина столбцов равна 8 (это позволяет оставить между столбцами небольшое расстояние, т. к. ширина интервалов равна 10). Сами столбцы смещены влево на 4 позиции, поэтому, например, левая и правая границы столбца "80" равны 76 и 84 соответственно, а его центр находится ровно на 80.

При помощи метода `plt.axis` мы задали отображение оси  $x$  в диапазоне от  $-5$  до 105 (чтобы столбцы "0" и "100" были полностью видны), а оси  $y$  — в диапазоне от 0 до 5. Метод `plt.xticks` устанавливает метки на оси  $x$  — 0, 10, 20, ..., 100.

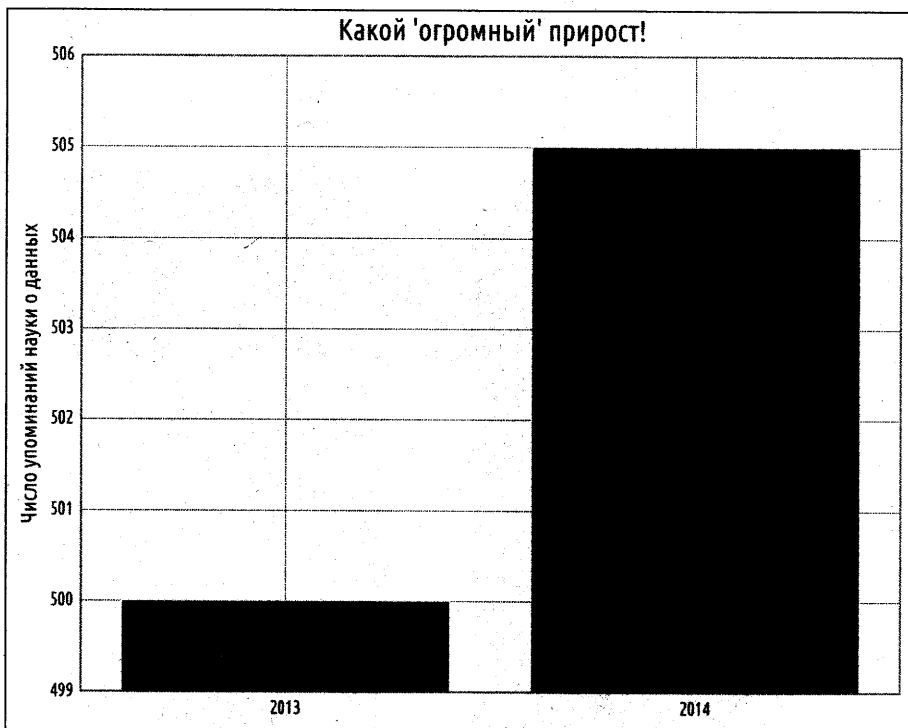


Рис. 3.4. Диаграмма с дезориентирующей осью  $y$

Будьте аккуратны при использовании метода `plt.axis`. Создание столбчатых диаграмм с осью  $y$ , которая начинается не с нуля, считается особенно дурным тоном, поскольку легко дезориентирует (рис. 3.4):

```
mentions = [500, 505]    # упоминания
years     = [2013, 2014] # годы
```

```
plt.bar([2012.6, 2013.6], mentions, 0.8)
plt.xticks(years)
plt.ylabel("Число упоминаний науки о данных")
```

```
# если этого не сделать, matplotlib подпишет ось X как 0, 1
# и добавит +2.013e3 в правом углу (недоработка matplotlib!)
plt.ticklabel_format(useOffset=False)
```

```
# дезориентирующая ось Y только показывает то, что выше 500
plt.axis([2012.5, 2014.5, 499, 506])
plt.title("Какой 'огромный' прирост!")
plt.show()
```

На рис. 3.5 использованы более осмысленные оси, и диаграмма выглядит уже не такой впечатляющей:

```
plt.axis([2012.5, 2014.5, 0, 550])
plt.title("Больше не такой огромный")
plt.show()
```

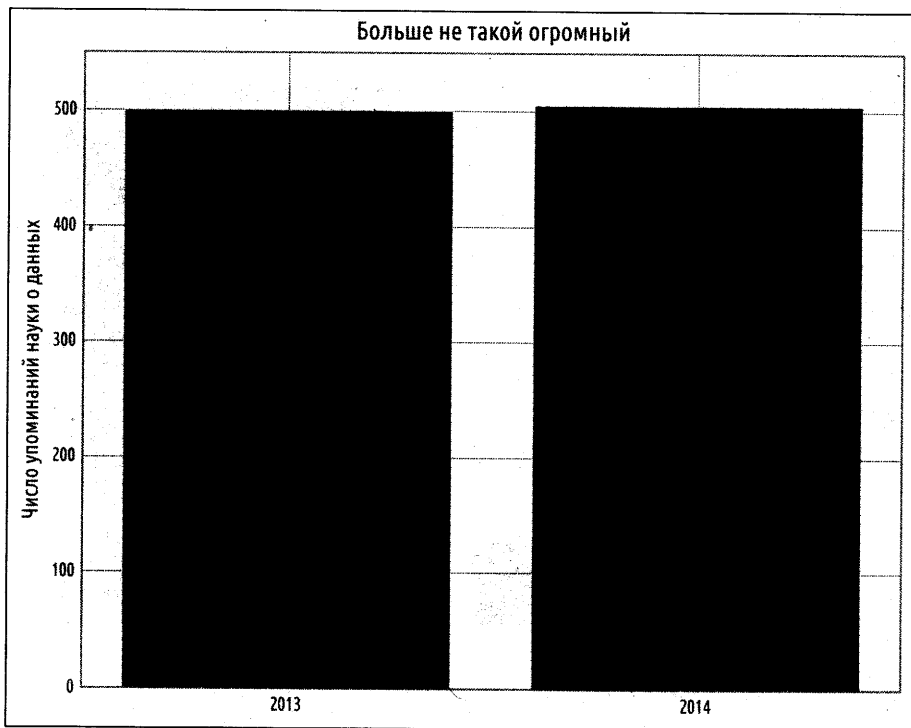


Рис. 3.5. Та же самая диаграмма с корректной осью y

## Линейные графики

Как мы уже убедились, линейные графики можно создавать при помощи метода `plt.plot`. Они хорошо подходят для изображения *тенденций* (рис. 3.6):

```
variance      = [1, 2, 4, 8, 16, 32, 64, 128, 256]      # дисперсия
bias_squared  = [256, 128, 64, 32, 16, 8, 4, 2, 1]     # квадрат смещения
# суммарная ошибка
total_error   = [x + y for x, y in zip(variance, bias_squared)]
xs = [i for i, _ in enumerate(variance)]
```

```
# метод plt.plot можно вызывать много раз,  
# чтобы показать несколько графиков на одной и той же диаграмме:  
# зеленая сплошная линия  
plt.plot(xs, variance, 'g-', label='дисперсия')  
# красная штрихпунктирная  
plt.plot(xs, bias_squared, 'r-.', label='смещение^2')  
# синяя пунктирная  
plt.plot(xs, total_error, 'b:', label='суммарная ошибка')  
  
# если для каждой линии задано название label,  
# то легенда будет показана автоматически,  
# loc=9 означает "наверху посередине"  
plt.legend(loc=9)  
plt.xlabel("Сложность модели")  
plt.title("Компромисс между смещением и дисперсией")  
plt.show()
```

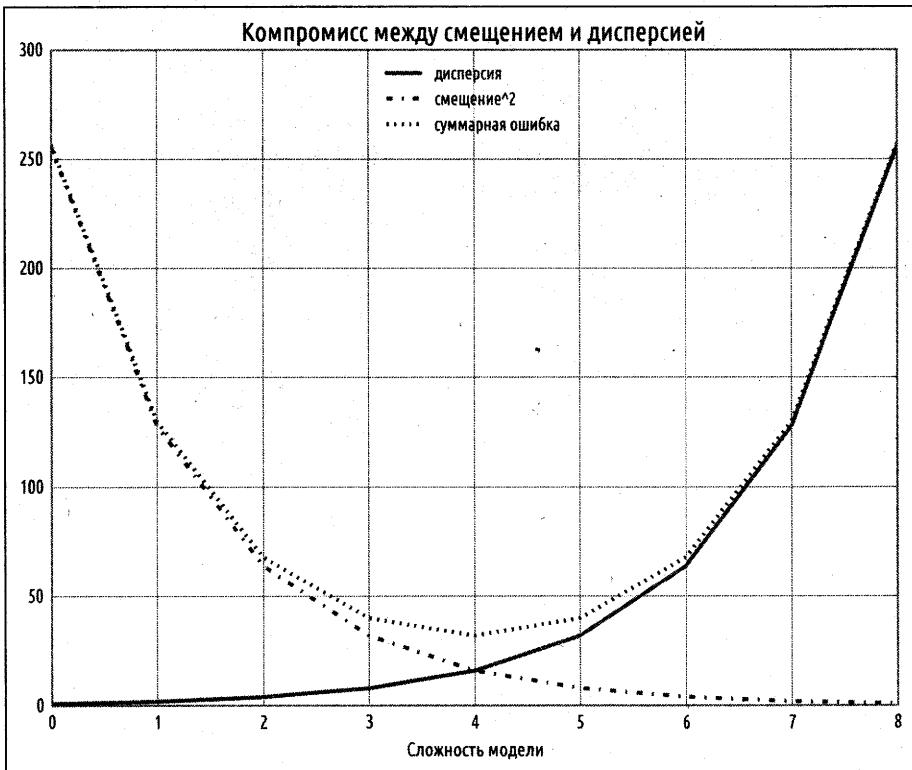


Рис. 3.6. Несколько линейных графиков с легендой

## Точечные диаграммы

Точечная диаграмма или *диаграмма рассеивания* (scatterplot) лучше всего подходит для визуализации связи между двумя парными выборками данных. Например, на рис. 3.7 показана связь между количеством друзей пользователя и количеством минут, которые он проводит на сайте каждый день:

```
friends = [ 70, 65, 72, 63, 71, 64, 60, 64, 67] # друзья
minutes = [175, 170, 205, 120, 220, 130, 105, 145, 190] # минуты
labels = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i'] # метки

plt.scatter(friends, minutes)

# назначить метку для каждой точки
for label, friend_count, minute_count in zip(labels, friends, minutes):
    plt.annotate(label,
                 xy=(friend_count, minute_count), # задать метку
                 xytext=(5, -5),                # и немного сместить ее
                 textcoords='offset points')

plt.title("Зависимость между количеством минут и числом друзей")
plt.xlabel("Число друзей")
plt.ylabel("Время, проводимое на сайте ежедневно, мин")
plt.show()
```

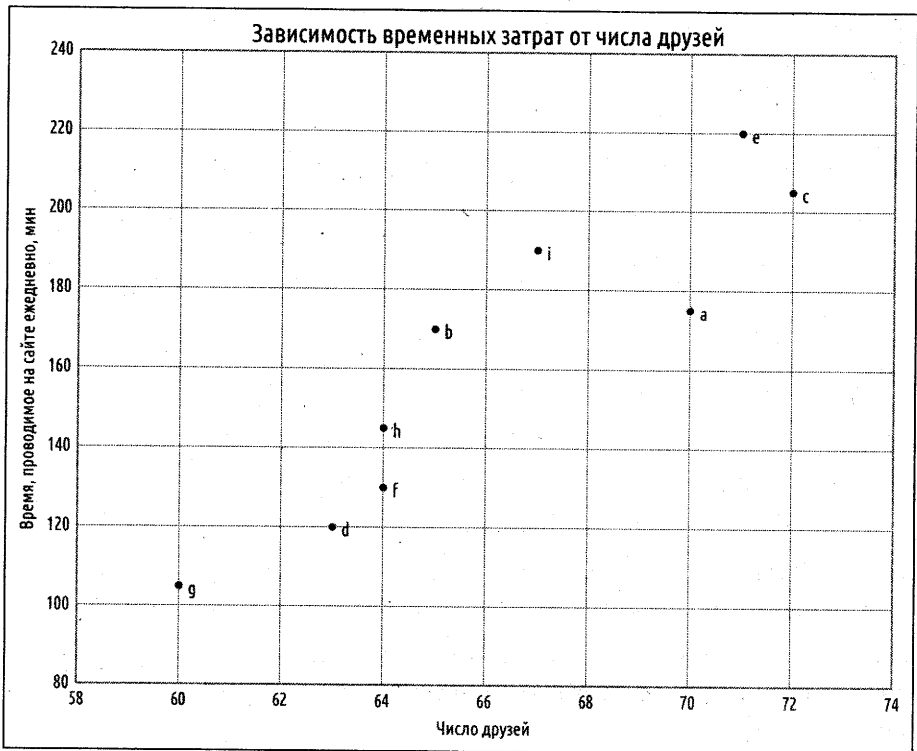


Рис. 3.7. Точечная диаграмма зависимости количества друзей от времени, проводимого на сайте

При размещении на точечной диаграмме сопоставимых переменных можно получить искаженную картину, если масштаб по осям будет определяться библиотекой matplotlib автоматически, как показано на рис. 3.8:

```
test_1_grades = [ 99, 90, 85, 97, 80] # оценки за тест 1
test_2_grades = [100, 85, 60, 90, 70] # оценки за тест 2

plt.scatter(test_1_grades, test_2_grades)
plt.title("Несопоставимые оси")
plt.xlabel("Оценки за тест № 1")
plt.ylabel("Оценки за тест № 2")
plt.show()
```

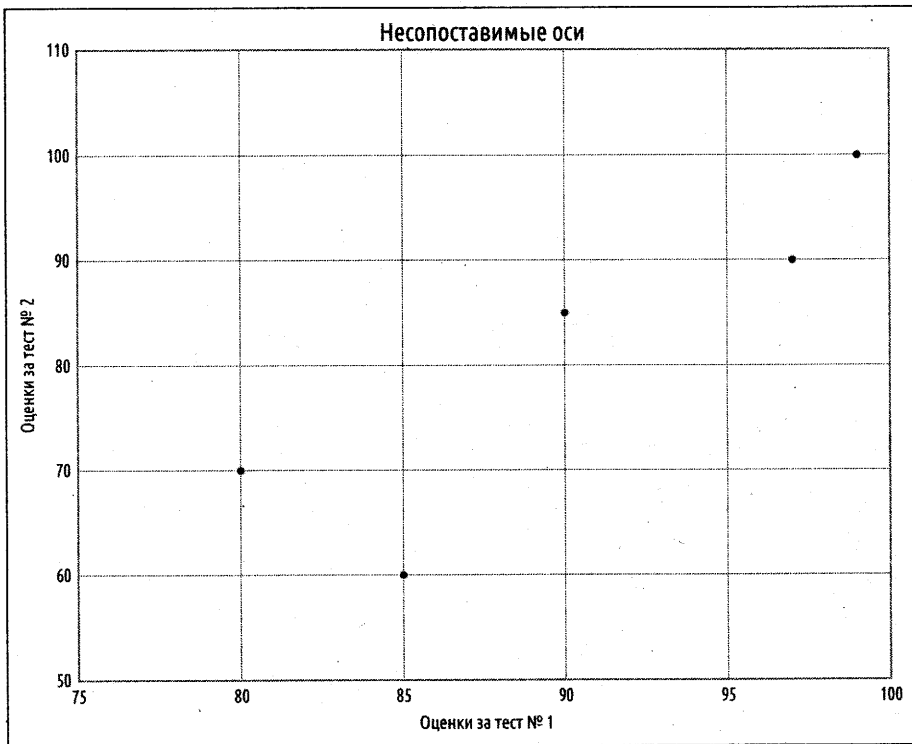


Рис. 3.8. Точечная диаграмма с несопоставимыми осями

Но если мы добавим вызов метода `plt.axis("equal")`, то диаграмма более точно покажет, что большая часть вариации оценок приходится на тест № 2 (рис. 3.9).

Этой информации вам будет достаточно для того, чтобы приступить к визуализации своих данных. Гораздо больше о методах визуализации мы узнаем в остальной части книги.

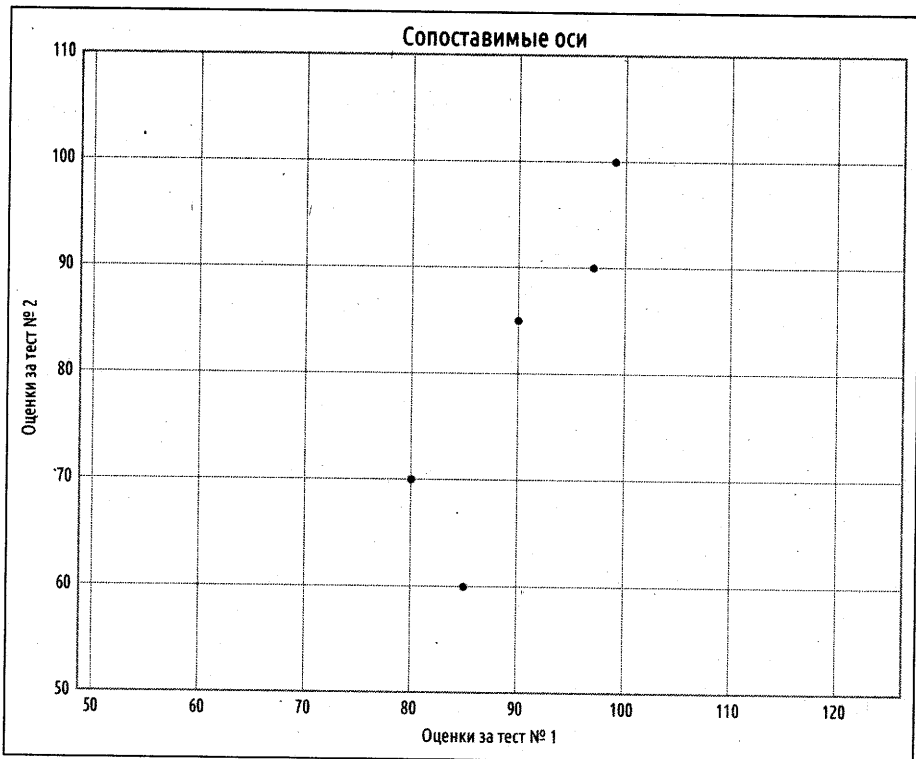


Рис. 3.9. Та же самая точечная диаграмма с сопоставимыми осями

## Для дальнейшего изучения

- ◆ Библиотека `seaborn` (<http://web.stanford.edu/~mwaskom/software/seaborn/>), разработанная на основе `matplotlib`, позволяет легко создавать более красивые и более сложные визуализации.
- ◆ Популярная библиотека `D3.js` на JavaScript (<https://d3js.org/>), предназначенная для создания продвинутых интерактивных визуализаций для Интернета, отражает последние тенденции в этой сфере. Несмотря на то, что она написана не на Python, вам стоит познакомиться с ней поближе.
- ◆ Библиотека `Bokeh` (<http://bokeh.pydata.org/>) — это сравнительно новая библиотека, которая позволяет создавать визуализации в стиле `D3.js` на Python.
- ◆ Библиотека `ggplot` (<https://pypi.python.org/pypi/ggplot>) — это портированная на Python версия популярной библиотеки `ggplot2`, написанной на языке R, которая широко применяется для создания графиков и диаграмм типографского качества. Если вы уже пользовались `ggplot2`, то библиотека `ggplot` будет, наверное, самой интересной для вас, а если нет, то, возможно, она покажется немного запутанной.



# Линейная алгебра

Есть ли что-то более бесполезное или менее полезное, чем алгебра?

Билли Коннолли<sup>1</sup>

*Линейная алгебра* — это область математики, которая изучает *векторные пространства*. В этой короткой главе перед вами не ставится задача разобраться с целым разделом алгебры, но так как линейная алгебра лежит в основе большинства понятий и методов науки о данных, я буду признателен, если вы хотя бы попытаетесь. Все, чему мы научимся в этой главе, очень пригодится нам в остальной части книги.

## Векторы

В абстрактном смысле *векторы* — это объекты, которые можно складывать между собой и умножать на *скалярные величины* (т. е. на числа), формируя новые векторы. Конкретно в нашем случае векторы — это точки в некотором конечномерном пространстве. Несмотря на то, что вы можете не рассматривать сами данные как векторы, ими удобно представлять последовательности чисел.

Например, если у вас есть данные о росте, весе и возрасте большого числа людей, то эти данные можно трактовать как трехмерные векторы (рост, вес, возраст). Если вы готовите группу учащихся к четырем экзаменам, то оценки студентов можно трактовать как четырехмерные векторы (экзамен1, экзамен2, экзамен3, экзамен4).

В самом простом случае векторы реализуют в виде списков чисел. Тогда список из трех чисел будет соответствовать вектору в трехмерном пространстве, и наоборот:

```
height_weight_age = [175, # сантиметры,
                     68,   # килограммы,
                     40 ] # годы
```

```
grades = [95, # экзамен1
          80, # экзамен2
          75, # экзамен3
          62 ] # экзамен4
```

---

<sup>1</sup> Билл Коннолли (1942) — шотландский комик, музыкант, ведущий и актер (см. [https://ru.wikipedia.org/wiki/Коннолли,\\_Билли](https://ru.wikipedia.org/wiki/Коннолли,_Билли)). — Прим. пер.

Одна из проблем с таким подходом заключается в том, что нам хотелось бы выполнять над ними *арифметические* операции, а поскольку списки в Python не являются векторами (и, следовательно, не располагают средствами векторной арифметики), то нам придется реализовывать эти арифметические операции самим. Поэтому начнем с них.

Нередко приходится *складывать два вектора*. Векторы складывают *покомпонентно*, т. е. если два вектора  $v$  и  $w$  имеют одинаковую длину, то их сумма — это вектор, чей первый элемент равен  $v[0] + w[0]$ , второй —  $v[1] + w[1]$  и т. д. (Если длины векторов разные, операция сложения не выполняется.)

Например, сложение векторов  $[1, 2]$  и  $[2, 1]$  даст  $[1 + 2, 2 + 1]$  или  $[3, 3]$ , как показано на рис. 4.1.

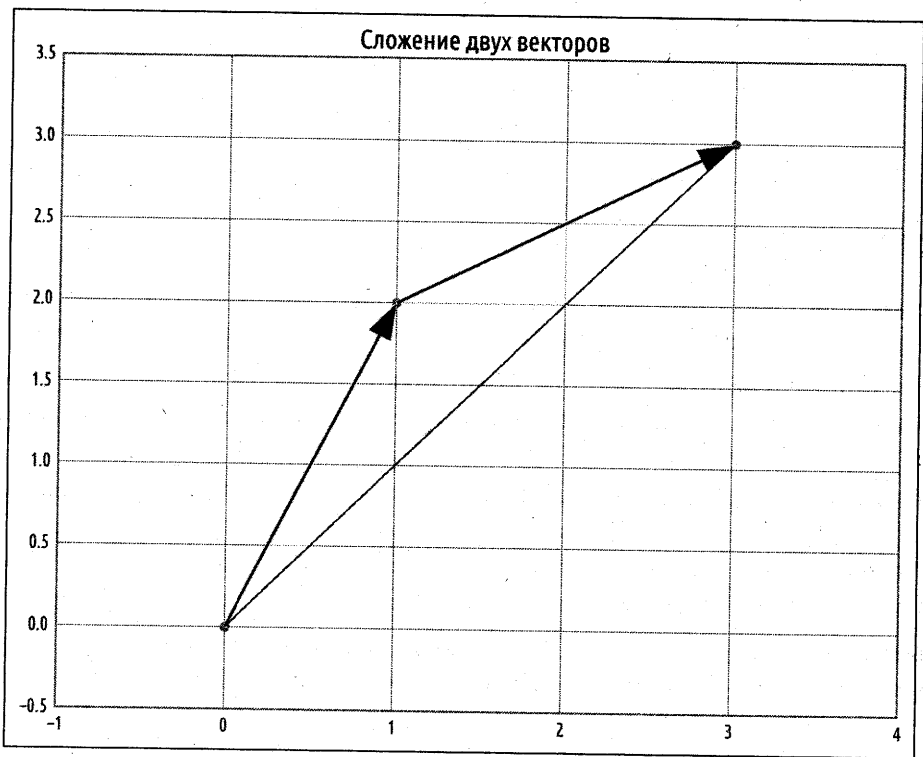


Рис. 4.1. Сложение двух векторов

Это легко реализуется объединением двух векторов посредством функции `zip` и сложением соответствующих элементов с помощью генератора последовательности:

```
def vector_add(v, w):
    """складывает соответствующие элементы"""
    return [v_i + w_i
            for v_i, w_i in zip(v, w)]
```

Вычитание двух векторов реализуется аналогичным образом, только берется разность соответствующих элементов:

```
def vector_subtract(v, w):
    """вычитает соответствующие элементы"""
    return [v_i - w_i
            for v_i, w_i in zip(v, w)]
```

Также может понадобиться *покомпонентная сумма списка векторов*, т. е. нужно создать новый вектор, чей первый элемент равен сумме всех первых элементов векторов, второй элемент — сумме всех вторых элементов и т. д. Простейший способ сделать это — сложить векторы один за другим:

```
def vector_sum(vectors):
    """суммирует все соответствующие элементы"""
    result = vectors[0] # начать с первого вектора
    for vector in vectors[1:]: # пройти в цикле по остальным векторам
        result = vector_add(result, vector) # и сложить их с результатом
    return result
```

Фактически, происходит свертка списка векторов при помощи `vector_add` (в качестве комбинирующей функции), поэтому можно написать короче, используя функции высшего порядка:

```
def vector_sum(vectors):
    return reduce(vector_add, vectors)
```

или даже так:

```
vector_sum = partial(reduce, vector_add)
```

Хотя последний способ, скорее, изощренный, чем пригодный для использования.

Кроме того, нам понадобится *умножение вектора на скаляр*, которое реализуется простым умножением каждого элемента вектора на это число:

```
def scalar_multiply(c, v):
    """c - это число, v - вектор"""
    return [c * v_i for v_i in v]
```

Это позволит вычислять покомпонентное среднее значение списка векторов (одинакового размера):

```
def vector_mean(vectors):
    """вычислить вектор, чей i-й элемент - это среднее значение
    всех i-х элементов входящих векторов"""
    n = len(vectors)
    return scalar_multiply(1/n, vector_sum(vectors))
```

Менее очевидная операция — это *скалярное произведение*. Скалярное произведение двух векторов есть сумма их покомпонентных произведений:

```
def dot(v, w):
    """v_1 * w_1 + ... + v_n * w_n"""
    return sum(v_i * w_i
               for v_i, w_i in zip(v, w))
```

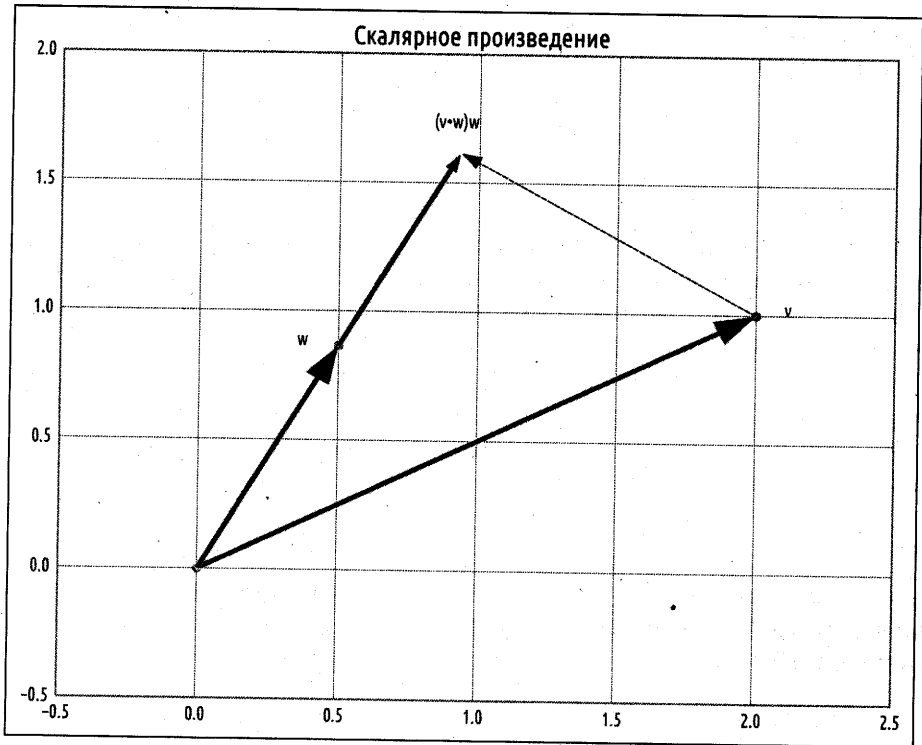


Рис. 4.2. Скалярное произведение как числовая проекция вектора

Скалярное произведение определяет величину, на которую вектор  $v$  простирается в направлении вектора  $w$ . Например, если  $w = [1, 0]$ , тогда скалярное произведение  $\text{dot}(v, w)$  — это просто первый компонент вектора  $v^2$ . Выражаясь иначе, это длина вектора, которая получится, если спроецировать вектор  $v$  на вектор  $w$  (рис. 4.2)<sup>3</sup>.

Используя эту функцию, легко вычислить сумму квадратов вектора:

```
def sum_of_squares(v):
    """v_1 * v_1 + ... + v_n * v_n"""
    return dot(v, v)
```

которую можно применить для вычисления его величины (или длины):

```
import math
```

```
def magnitude(v):
    return math.sqrt(sum_of_squares(v)) # math.sqrt - квадратный корень
```

<sup>2</sup> Компонент вектора — это величина проекции на базис, где базис — набор векторов, образующих оси координат. Любой вектор можно выразить через проекции его на базисные векторы, эти проекции — компоненты вектора. Например, если задан некоторый базис  $i, j, k$ , то вектор  $\vec{x}$  представляется в виде:  $\vec{x} = ai + bj + ck$ . Здесь  $a, b, c$  и есть компоненты вектора  $\vec{x}$ . — Прим. пер.

<sup>3</sup> Для справки: в координатной форме скалярным произведением двух векторов  $v$  и  $w$  называется произведение длины вектора  $w$  на числовую проекцию вектора  $v$  на направление вектора  $w$  или произведение длины вектора  $v$  на числовую проекцию вектора  $w$  на направление вектора  $v$ . — Прим. пер.

Теперь у нас есть все функции, необходимые для вычисления расстояния (точнее, евклидова расстояния) между двумя векторами по формуле:

$$\sqrt{(v_1 - w_1)^2 + \dots + (v_n - w_n)^2}.$$

# квадрат расстояния между двумя векторами

```
def squared_distance(v, w):
    """(v_1 - w_1) ** 2 + ... + (v_n - w_n) ** 2"""
    return sum_of_squares(vector_subtract(v, w))
```

# расстояние между двумя векторами

```
def distance(v, w):
    return math.sqrt(squared_distance(v, w))
```

Вероятно, будет понятнее, если мы запишем эту функцию в таком виде (что вычисляет то же самое):

```
def distance(v, w):
    return magnitude(vector_subtract(v, w))
```

Для начала этого вполне достаточно. Мы будем пользоваться этими функциями на протяжении всей книги.



Использование списков как векторов великолепно подходит в качестве иллюстрации, но крайне неэффективно по производительности.

В рабочем коде следует воспользоваться библиотекой NumPy, в которой есть класс, реализующий высокоэффективный массив с разными арифметическими операциями.

## Матрицы

*Матрица* — это двумерный ряд чисел. Мы будем реализовывать матрицы как списки списков, где все внутренние списки имеют одинаковый размер и представляют собой *строку* матрицы. Если  $A$  — это матрица, то  $A[i][j]$  — это элемент в  $i$ -й строке и  $j$ -м столбце. В соответствии с правилами математической записи для обозначения матриц мы будем, в основном, использовать заглавные буквы. Например:

```
A = [[1, 2, 3],      # матрица A имеет 2 строки и 3 столбца
      [4, 5, 6]]
```

```
B = [[1, 2],        # матрица B имеет 3 строки и 2 столбца
      [3, 4],
      [5, 6]]
```



В математике, как правило, строки и столбцы нумеруются начиная с единицы ("первая строка", "первый столбец"). Но поскольку мы реализуем матрицы в виде списков Python, которые являются нуль-индексными, то будем называть первую строку матрицы "нулевой строкой", а первый столбец — "нулевым столбцом".

Представленная в виде списка списков матрица  $A$  имеет  $\text{len}(A)$  строк и  $\text{len}(A[0])$  столбцов. Эти значения образуют *форму* матрицы:

```
def shape(A):
    num_rows = len(A)
    num_cols = len(A[0]) if A else 0 # число элементов в первой строке
    return num_rows, num_cols
```

Если матрица имеет  $n$  строк и  $k$  столбцов, то мы будем говорить, что это матрица размера  $n \times k$ . Каждая строка матрицы размера  $n \times k$  может быть представлена вектором длиной  $k$ , а каждый столбец — вектором длиной  $n$ :

```
# получить i-ю строку
def get_row(A, i):
    return A[i] # A[i] - это i-я строка
```

```
# получить j-й столбец
def get_column(A, j):
    return [A_i[j] # j-й элемент строки A_i
            for A_i in A] # для каждой строки A_i
```

Кроме того, нам потребуется возможность создавать матрицу при наличии ее формы и функции, которая генерирует ее элементы. Это можно сделать на основе вложенного генератора последовательности:

```
def make_matrix(num_rows, num_cols, entry_fn):
    """возвращает матрицу размером num_rows x num_cols,
    (i,j)-й элемент которой равен функции entry_fn(i, j)"""
    return [[entry_fn(i, j) # при заданном i создать список
              for j in range(num_cols)] # [entry_fn(i,0), ... ]
            for i in range(num_rows)] # повторить для каждого i
```

При помощи этой функции можно, например, создать *единичную матрицу* размера  $5 \times 5$  (с единицами по диагонали и нулями в остальных элементах):

```
def is_diagonal(i, j):
    """единицы по диагонали, остальные нули"""
    return 1 if i == j else 0
```

```
# единичная матрица
identity_matrix = make_matrix(5, 5, is_diagonal)
```

```
# [[1, 0, 0, 0, 0],
# [0, 1, 0, 0, 0],
# [0, 0, 1, 0, 0],
# [0, 0, 0, 1, 0],
# [0, 0, 0, 0, 1]]
```

Матрицы имеют особое значение для нас по нескольким причинам.

Во-первых, мы можем использовать матрицу для представления набора данных, состоящего из нескольких векторов, рассматривая каждую строку матрицы в каче-

стве вектора. Например, если имеются данные о росте, весе и возрасте 1000 человек, то их можно представить в виде матрицы размера  $1000 \times 3$ :

```
data = [[175, 68, 40],
        [163, 64, 26],
        [193, 78, 19],
        # ...
    ]
```

Во-вторых, как мы убедимся далее, матрицу размера  $n \times k$  можно использовать в качестве линейной функции, которая отображает  $k$ -мерные векторы в  $n$ -мерные. Некоторые наши методы и концепции при анализе данных будут использовать такие функции.

И, в-третьих, матрицы можно использовать для двоичного представления дружеских связей. В главе 1 дуги социальной сети были представлены в виде множества пар  $(i, j)$ . Альтернативная реализация подразумевает создание матрицы смежности, т. е. матрицы  $A$ , такой, что элемент  $A[i][j]$  равен 1, если узлы  $i$  и  $j$  связаны между собой, и 0 — в противном случае.

Вспомним, что ранее мы представляли дружеские связи в виде списка кортежей:

```
friendships = [(0, 1), (0, 2), (1, 2), (1, 3), (2, 3), (3, 4),
               (4, 5), (5, 6), (5, 7), (6, 8), (7, 8), (8, 9)]
```

То же самое можно реализовать следующим образом:

```
# пользователь 0 1 2 3 4 5 6 7 8 9
#
friendships = [[0, 1, 1, 0, 0, 0, 0, 0, 0, 0], # пользователь 0
               [1, 0, 1, 1, 0, 0, 0, 0, 0, 0], # пользователь 1
               [1, 1, 0, 1, 0, 0, 0, 0, 0, 0], # пользователь 2
               [0, 1, 1, 0, 1, 0, 0, 0, 0, 0], # пользователь 3
               [0, 0, 0, 1, 0, 1, 0, 0, 0, 0], # пользователь 4
               [0, 0, 0, 0, 1, 0, 1, 1, 0, 0], # пользователь 5
               [0, 0, 0, 0, 0, 1, 0, 0, 1, 0], # пользователь 6
               [0, 0, 0, 0, 0, 1, 0, 0, 1, 0], # пользователь 7
               [0, 0, 0, 0, 0, 0, 1, 1, 0, 1], # пользователь 8
               [0, 0, 0, 0, 0, 0, 0, 0, 1, 0]] # пользователь 9
```

При очень небольшом количестве дружеских связей такое представление имеет гораздо меньшую эффективность, т. к. приходится хранить много нулей. Однако в матричном представлении проверка связи между двумя узлами выполняется намного быстрее — вместо того, чтобы проверять каждую дугу в списке, нужно всего лишь выполнить поиск в матрице:

```
friendships[0][2] == 1 # True, 0 и 2 - друзья
friendships[0][8] == 1 # False, 0 и 8 - не друзья
```

Аналогичным образом, чтобы найти все связи для конкретного узла, вам нужно лишь проверить соответствующий этому узлу столбец (или строку) матрицы:

```

friends_of_five = [i                                     # нужно обратиться
                    for i, is_friend in enumerate(friendships[5])
                    if is_friend]                       # лишь к одной строке

```

Ранее, чтобы ускорить этот процесс, мы добавляли к каждому узловому объекту список связей. Но для крупного эволюционирующего графа, представленного таким образом, это было бы слишком ресурсоемко и сложно.

Мы будем возвращаться к теме матриц на протяжении всей книги.

## Для дальнейшего изучения

- ◆ Линейная алгебра активно используется аналитиками данных (часто неявным образом и не вполне разбирающимися в ней). Поэтому неплохо было бы прочитать учебники по этой тематике. В Интернете можно найти бесплатные пособия, такие как:
  - "Линейная алгебра" от Калифорнийского университета в Дэвисе (<https://www.math.ucdavis.edu/~linear/>);
  - "Линейная алгебра" от Вермонтского колледжа в Сан-Мишель (<http://joshua.smcvt.edu/linearalgebra/>);
  - более продвинутое введение "Неправильная линейная алгебра" (<http://www.math.brown.edu/~treil/papers/LADW/LADW.html>) для тех, кто не боится сложностей.
- ◆ Весь функционал, который мы разработали в этой главе, а также очень много дополнительных возможностей можно найти в бесплатной библиотеке NumPy (<http://www.numpy.org/>).



# Статистика

Факты — упрямая вещь, но статистика гораздо стоворчивее.

Марк Твен<sup>1</sup>

*Статистика* опирается на математику и методы, при помощи которых мы пытаемся понять данные. Она представляет собой богатую и обширную область знаний, для которой больше подошла бы книжная полка или целая комната в библиотеке, а не глава в книге, поэтому наше изложение, конечно же, будет кратким. Тем не менее, представленного в этой главе материала вам будет достаточно, чтобы стать опасным человеком и пробудить у вас интерес для самостоятельного изучения.

## Описание одиночного набора данных

Благодаря полезному сочетанию живого слова и удачи, социальная сеть DataSciencester выросла до нескольких десятков пользователей, и директор по привлечению финансовых ресурсов просит вас проанализировать, сколько друзей есть у пользователей сети, чтобы он мог включить эти данные в свои "презентации для лифта"<sup>2</sup>.

Используя простые методы из *главы 1*, вы легко можете предъявить запрашиваемые данные. Однако сейчас вы столкнулись с задачей выполнения их *описательного анализа*.

Любой набор данных очевидным образом характеризует сам себя:

```
# число друзей
num_friends = [100, 49, 41, 40, 25,
               # ... и еще много других
               ]
```

И такое описание может оказаться наилучшим для сравнительного небольшого набора данных. Но для более объемного набора данных это будет выглядеть очень громоздко и, скорее всего, непрозрачно. (Представьте, что у вас перед глазами список из 1 млн чисел.) По этой причине пользуются статистиками (статистическими

---

<sup>1</sup> Первая часть афоризма была популяризована Марком Твеном, вторая приписывается Лоренсу Питеру (1919–1990) — канадско-американскому педагогу и литератору (см. [http://dic.academic.ru/dic.nsf/dic\\_wingwords/3596/Факты](http://dic.academic.ru/dic.nsf/dic_wingwords/3596/Факты)). — *Прим. пер.*

<sup>2</sup> Презентация для лифта — короткий рассказ о концепции продукта, проекта или сервиса, который можно полностью представить во время поездки на лифте. Используется для изложения концепции нового бизнеса в целях получения инвестиций. — *Прим. пер.*

показателями), при помощи которых обобщают и передают информацию о существенных признаках, присутствующих в данных.

Вначале вы помещаете количество друзей на гистограмму, используя словарь Counter и метод plt.bar (рис. 5.1):

```
friend_counts = Counter(num_friends)
xs = range(101)                                     # максимальное значение 100
ys = [friend_counts[x] for x in xs]                 # высота - это количество друзей
plt.bar(xs, ys)
plt.axis([0, 101, 0, 25])
plt.title("Гистограмма количества друзей")
plt.xlabel("Количество друзей")
plt.ylabel("Количество людей")
plt.show()
```

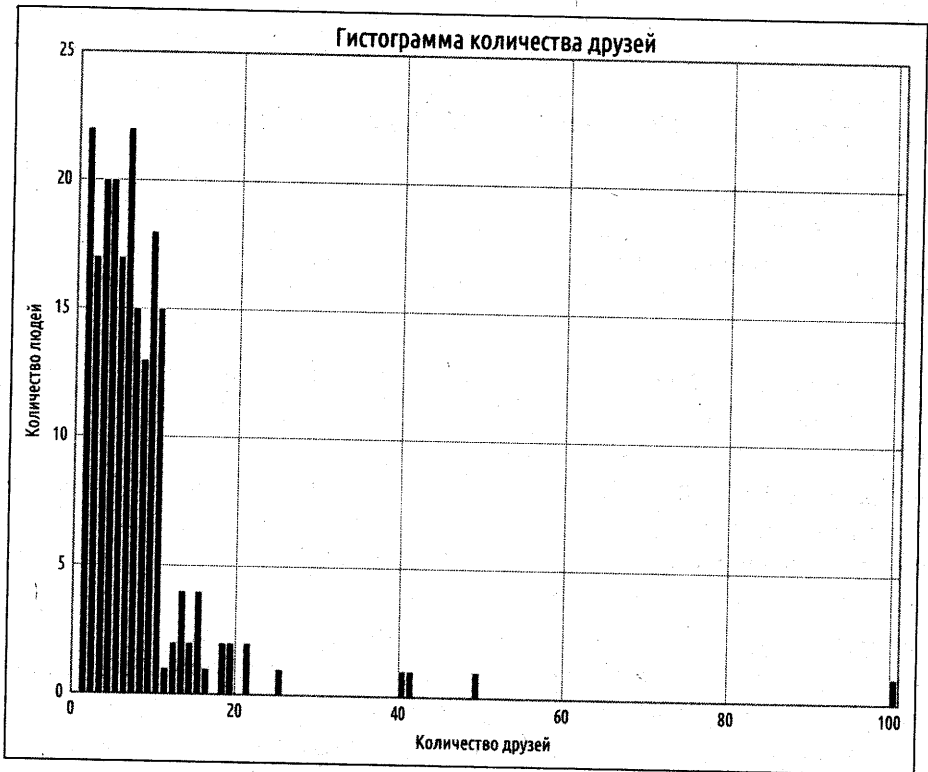


Рис. 5.1. Гистограмма количества друзей

К сожалению, эта диаграмма мало информативна, поэтому вы приступаете к формированию некоторых статистик. Наверное, самой простой статистикой является число точек данных:

```
num_points = len(num_friends)                       # число точек = 204
```

Кроме этого, могут быть интересны наибольшие и наименьшие значения:

```
largest_value = max(num_friends) # максимальное значение = 100
smallest_value = min(num_friends) # минимальное значение = 1
```

которые являются особыми случаями, когда нужно узнать значения в определенных позициях:

```
sorted_values = sorted(num_friends) # отсортированные значения
smallest_value = sorted_values[0] # минимум = 1
second_smallest_value = sorted_values[1] # следующий минимум = 1
second_largest_value = sorted_values[-2] # следующий максимум = 49
```

Но это только начало.

## Показатели центра распределения

Обычно мы хотим иметь некоторое представление о том, где находится центр данных. Чаще всего для этих целей используется *среднее* (или среднее арифметическое) значение, которое берется как сумма данных, деленная на их количество<sup>3</sup>:

```
# среднее значение
# В Python 2 частное округляется вниз, поэтому следует импортировать
# вещественное деление from __future__ import division
def mean(x):
    return sum(x) / len(x)
```

```
mean(num_friends) # 7.333333
```

Для двух точек средней является точка, лежащая посередине между ними. По мере добавления других точек среднее значение будет смещаться в разные стороны, в зависимости от значения каждой новой точки.

Кроме среднего значения, иногда может понадобиться *медиана*, которая является ближайшим к центру значением (если число точек данных нечетное) либо средним арифметическим, взятым как полусумма двух ближайших к центру значений (если число точек четное).

Например, если у нас есть отсортированный вектор  $x$  из пяти точек, то медианой будет  $x[5 // 2]$  или  $x[2]$ . Если в векторе шесть точек, то берется среднее арифметическое  $x[2]$  (третья точка) и  $x[3]$  (четвертая точка).

Обратите внимание, что медиана — в отличие от среднего — не зависит от каждого значения в наборе данных. Например, если сделать наибольшую точку еще больше (или наименьшую точку еще меньше), то срединные точки останутся неизменными, следовательно, и медиана не изменится.

<sup>3</sup> Более общий вид имеет формула средней арифметической *взвешенной*, где каждое значение переменной сначала умножается на свой вес. В то время как средняя арифметическая *простая* — это частный случай, когда веса для всех значений равны 1. — Прим. пер.

Функция `median` имеет более сложную реализацию, чем можно было бы ожидать, в основном из-за того, что приходится учитывать случай с четностью:

```
def median(v):
    """возвращает ближайшее к середине значение для v"""
    n = len(v)
    sorted_v = sorted(v)
    midpoint = n // 2 # индекс срединного значения

    if n % 2 == 1:
        # если нечетное, вернуть срединное значение
        return sorted_v[midpoint]
    else:
        # если четное, вернуть среднее 2-х срединных значений
        lo = midpoint - 1
        hi = midpoint
        return (sorted_v[lo] + sorted_v[hi]) / 2

median(num_friends) # = 6.0
```

Среднее значение, конечно, вычисляется проще, и оно несколько варьирует по мере изменения данных. Если у нас есть  $n$  точек, и одна из них увеличилась на любое малое число  $e$ , то среднее обязательно увеличится на  $e/n$ . (Этот факт делает его подверженным разного рода ухищрениям при калькуляции.) А для того чтобы найти медиану, данные нужно сперва отсортировать, и изменение одной из точек на любое малое число  $e$  может увеличить медиану на величину равную  $e$ , меньшую чем  $e$ , либо не изменить совсем (в зависимости от всего набора данных).



На самом деле существуют неочевидные приемы эффективного вычисления медиан<sup>4</sup> без сортировки данных. Поскольку их рассмотрение выходит за рамки этой книги, мы вынуждены будем сортировать данные.

Вместе с тем, среднее значение очень чувствительно к *выбросам* в данных. Если бы самый дружелюбный пользователь имел 200 друзей (вместо 100), то среднее увеличилось бы до 7.82, а медиана осталась бы на прежнем уровне. Так как выбросы являются, скорее всего, плохими данными (или иначе — нерепрезентативными для ситуации, в которой мы пытаемся разобраться), то среднее может иногда давать искаженную картину. В качестве примера часто приводят историю, как в середине 1980-х годов в Университете Северной Каролины география стала специализацией с самой высокой стартовой среднестатистической зарплатой, что произошло в основном из-за звезды НБА (и выброса) Майкла Джордана.

Обобщением медианы является *квантиль* — значение, меньше которого расположен определенный процент<sup>5</sup> данных. (Медиана — это значение, меньше которого расположены 50% данных.)

<sup>4</sup> <https://en.wikipedia.org/wiki/Quickselect>.

<sup>5</sup> Под процентом здесь на самом деле подразумевается *процентиль*, т. е. значение в упорядоченной выборке, ниже которого расположен заданный процент данных. — Прим. пер.

```
# квантиль
def quantile(x, p):
    """возвращает значение в x, соответствующее p-ому проценту данных"""
    p_index = int(p * len(x)) # преобразует значение % в индекс списка
    return sorted(x)[p_index]

quantile(num_friends, 0.10) # 1
quantile(num_friends, 0.25) # 3 (нижний квантиль)
quantile(num_friends, 0.75) # 9 (верхний квантиль)
quantile(num_friends, 0.90) # 13
```

Реже используют *моду* — значение или значения, которые встречаются наиболее часто:

```
# мода
def mode(x):
    """возвращает список, поскольку мод может быть больше одной"""
    counts = Counter(x)
    max_count = max(counts.values())
    return [x_i for x_i, count in counts.iteritems()
            if count == max_count]

mode(num_friends) # 1 и 6
```

Но наиболее часто мы будем использовать просто среднее значение.

## Показатели вариации

Показатели *вариации* отражают меру изменчивости данных. Как правило, это статистические показатели, для которых значения, близкие к нулю, означают полное *отсутствие изменчивости*, а большие значения (что бы это ни означало) — *очень большую изменчивость*. Например, самым простым показателем является *размах*, который определяется как разница между максимальным и минимальным значениями данных:

```
# размах
# слово "range" в Python уже используется, поэтому берем другое
def data_range(x):
    return max(x) - min(x)

data_range(num_friends) # 99
```

Размах равен нулю, когда  $\max$  и  $\min$  одинаковые, что происходит только тогда, когда все элементы  $x$  равны между собой, и значит, изменчивость в данных отсутствует. И наоборот, когда размах широкий, то  $\max$  много больше  $\min$ , и изменчивость в данных высокая.

Как и медиана, размах также не зависит от всего набора данных. Набор данных, все точки которого равны 0 или 100, имеет тот же размах, что и набор данных, чьи зна-

чения представлены числами 0, 100 и серией из числа 50, хотя кажется, что первый набор "должен" варьировать больше.

Более точным показателем вариации является *дисперсия*, вычисляемая как<sup>6</sup>:

```
# вектор отклонений от среднего (центрировать вектор)
def de_mean(x):
    """пересчитать x, вычтя его среднее (среднее результата будет = 0)"""
    x_bar = mean(x) # среднее
    return [x_i - x_bar for x_i in x]

# дисперсия - это средняя сумма квадратов отклонений от среднего;
# в зарубежной литературе этот показатель носит название variance
def variance(x):
    """предполагается, что вектор x состоит как минимум из 2-х элементов"""
    n = len(x)
    deviations = de_mean(x)
    return sum_of_squares(deviations) / (n - 1)

variance(num_friends) # 81.54
```



Выглядит почти как средний квадрат отклонений от средних значений, но с одним исключением: в знаменателе не  $n$ , а на  $n - 1$ . В действительности, когда имеют дело с выборкой из более крупной генеральной совокупности, переменная  $x\_bar$  является лишь *приблизительной оценкой* среднего, где  $(x\_i - x\_bar) ** 2$  в среднем дает заниженную оценку квадрата отклонения от среднего для  $x\_i$ . Поэтому делят не на  $n$ , а на  $(n - 1)$ <sup>7</sup>. См. Википедию<sup>8</sup>.

При этом, в каких бы единицах ни измерялись данные (в "друзьях", например), все показатели центра распределения вычисляются в тех же самых единицах измерения. Аналогичная ситуация и с размахом. Дисперсия же измеряется в единицах, которые представляют собой *квадрат* исходных единиц ("друзья в квадрате"). Поскольку такие единицы измерения трудно интерпретировать, то вместо дисперсии мы будем чаще обращаться к *стандартному отклонению* (корень из дисперсии):

```
# стандартное отклонение
def standard_deviation(x):
    return math.sqrt(variance(x))

standard_deviation(num_friends) # 9.03
```

<sup>6</sup> Функция `variance` вычисляет несмещенную оценку выборочной дисперсии ( $n - 1$  в знаменателе). — Прим. пер.

<sup>7</sup> Поскольку соотношение между выборочной и генеральной дисперсией составляет  $n/(n - 1)$ , то с ростом  $n$  данное выражение стремится к 1, т. е. разница между значениями выборочной и генеральной дисперсиями уменьшается. На практике при  $n > 30$  уже нет разницы, какое число стоит в знаменателе:  $n$  или  $n - 1$ . — Прим. пер.

<sup>8</sup> [https://en.wikipedia.org/wiki/Unbiased\\_estimation\\_of\\_standard\\_deviation](https://en.wikipedia.org/wiki/Unbiased_estimation_of_standard_deviation).

Размах и стандартное отклонение так же чувствительны к выбросам, как и среднее. На том же самом примере, если бы у самого дружелюбного пользователя было 200 друзей, то стандартное отклонение было бы 14.89 или на 60% больше!

Более надежной альтернативой является вычисление *интерквартильного размаха* или разности между значением, соответствующим 75%, и значением, соответствующим 25% данных:

```
# интерквартильный размах
def interquartile_range(x):
    return quantile(x, 0.75) - quantile(x, 0.25)
```

```
interquartile_range(num_friends) # 6
```

Этот показатель позволяет простым образом исключить влияние небольшого числа выбросов.

## Корреляция

У директора по развитию сети есть теория, что количество времени, которое люди проводят на сайте, связано с числом их друзей (она же не просто так занимает должность директора), и она просит вас проверить это предположение.

Изучив журналы трафика, вы создали список `daily_minutes`, который показывает, сколько минут в день каждый пользователь тратит на DataSciencester, и упорядочили его так, чтобы его элементы соответствовали элементам нашего предыдущего списка `num_friends`. Нам нужно найти связь между этими двумя метриками.

Сперва обратимся к *ковариации* — парному аналогу дисперсии. В отличие от дисперсии, которая измеряет отклонение одной переменной от ее среднего, ковариация измеряет совместное отклонение двух переменных от своих средних:

```
# ковариация def covariance(x, y):
    n = len(x)
    return dot(de_mean(x), de_mean(y)) / (n - 1)
```

```
covariance(num_friends, daily_minutes) # 22.43
```

Напомним, что функция `dot` суммирует произведения соответствующих пар элементов (см. главу 3). Когда соответствующие элементы векторов  $x$  и  $y$  оба одновременно выше или ниже своих средних, то в сумму входит положительное число. Когда один из них находится выше своего среднего, а другой — ниже, то в сумму входит отрицательное число. Следовательно, "большая" положительная ковариация означает, что  $x$  стремится принимать большие значения при больших значениях  $y$  и малые значения — при малых значениях  $y$ . "Большая" отрицательная ковариация означает обратное —  $x$  стремится принимать малые значения при большом  $y$  и наоборот. Ковариация, близкая к нулю, означает, что такой связи не существует.

Тем не менее, этот показатель бывает трудно интерпретировать, и вот почему:

- ◆ единицами измерения ковариации являются произведения единиц входящих переменных (например, число друзей и минуты в день), которые трудно понять (что такое "друг в минуту в день?");
- ◆ если бы у каждого пользователя было в 2 раза больше друзей (но такое же количество минут, проведенных на сайте), то ковариация была бы в 2 раза больше. Однако в некотором смысле степень взаимосвязи между ними осталась бы на прежнем уровне. Говоря иначе, трудно определить, что считать "большой" ковариацией.

Поэтому чаще обращаются к *корреляции*, в которой ковариация распределяется между стандартными отклонениями обеих переменных:

```
# корреляция
def correlation(x, y):
    stdev_x = standard_deviation(x)
    stdev_y = standard_deviation(y)
    if stdev_x > 0 and stdev_y > 0:
        return covariance(x, y) / stdev_x / stdev_y
    else:
        return 0    # если переменные не меняются, то корреляция равна нулю
```

```
correlation(num_friends, daily_minutes)    # 0.25
```

Корреляция является безразмерной величиной, ее значения всегда лежат между  $-1$  (идеальная антикорреляция) и  $1$  (идеальная корреляция). Так, число  $0.25$  представляет собой относительно слабую положительную корреляцию.

Впрочем, мы забыли сделать одну вещь — проверить наши данные. Посмотрим на рис. 5.2.

Человек, у которого 100 друзей, проводит на сайте всего одну минуту в день, и поэтому он представляет собой "дикий" выброс, а корреляция, как будет показано далее, может быть очень чувствительна к выбросам. Что произойдет, если мы проигнорируем этот выброс?

```
outlier = num_friends.index(100)    # индекс выброса

# отфильтровать выброс
num_friends_good = [x
                    for i, x in enumerate(num_friends)
                    if i != outlier]

daily_minutes_good = [x
                    for i, x in enumerate(daily_minutes)
                    if i != outlier]

correlation(num_friends_good, daily_minutes_good)    # 0.57
```

Без выброса получается более сильная корреляция (рис. 5.3).



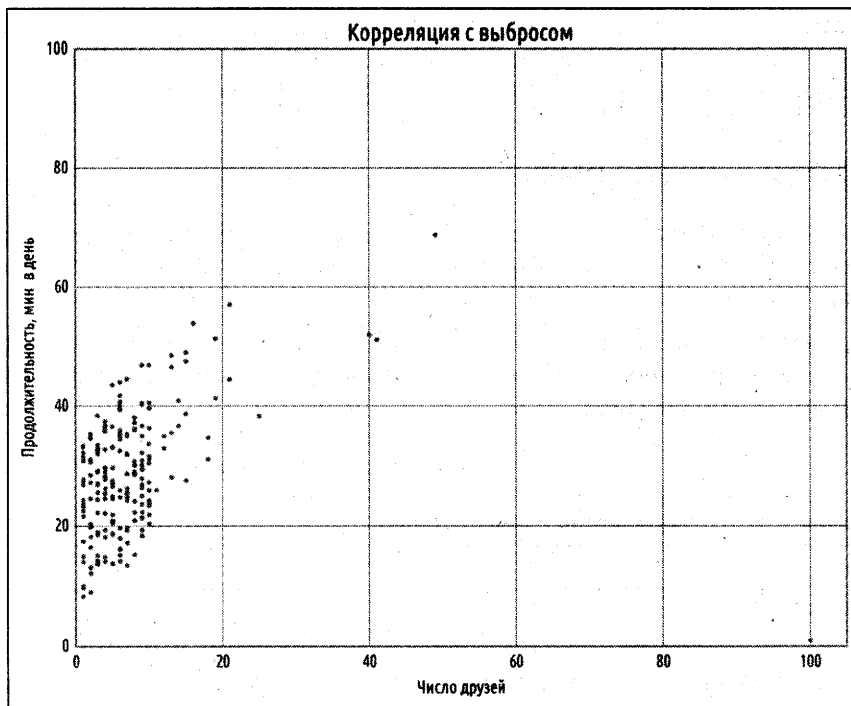


Рис. 5.2. Корреляция с выбросом

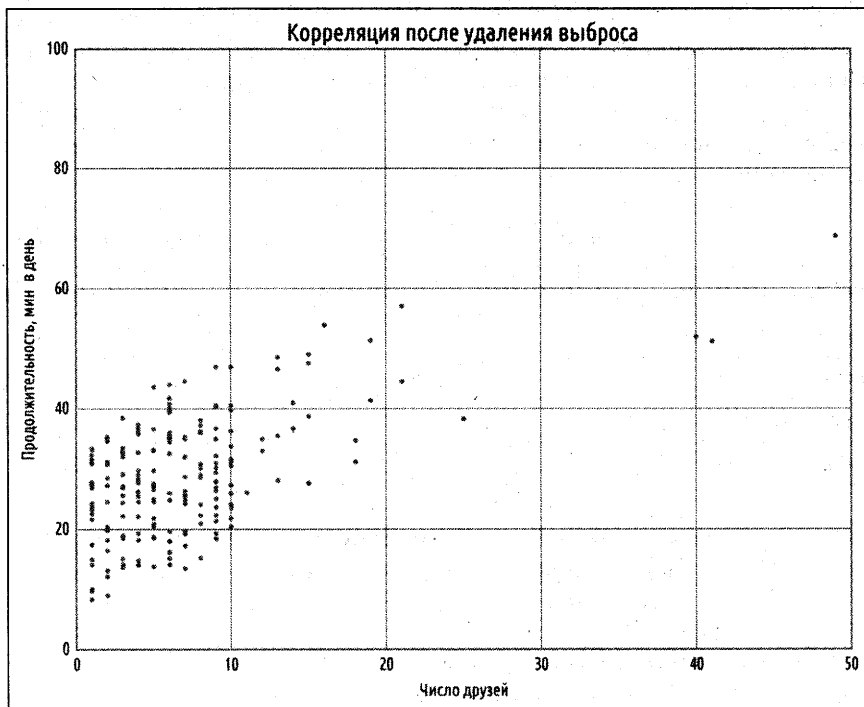


Рис. 5.3. Корреляция после удаления выброса

В результате дальнейших исследований вы обнаруживаете, что выброс на самом деле был техническим *пробным* аккаунтом, который забыли удалить, и вы поступили совершенно правильно, исключив его из рассмотрения.

## Парадокс Симпсона

Нередко при анализе данных вызывает удивление парадокс Симпсона<sup>9</sup>, при котором корреляции могут быть обманчивыми, если игнорируются *спутывающие* переменные<sup>10</sup>.

Например, вы разделили всех пользователей сети на аналитиков с Восточного побережья или с Западного побережья и решаете выяснить, с какого побережья аналитики дружелюбнее (табл. 5.1).

Таблица 5.1. География пользователей сети

Побережье	Количество пользователей	Среднее число друзей
Западное	101	8,2
Восточное	103	6,5

Очевидно, что аналитики с Западного побережья дружелюбнее, чем аналитики с Восточного побережья. Ваши коллеги выдвигают разного рода предположения, почему так происходит: может быть, дело в солнце, или в кофе, или в органических продуктах, или в безмятежной атмосфере Тихого океана?

Продолжая изучать данные, вы обнаруживаете что-то очень необычное. Если учитывать только пользователей с ученой степенью, то у аналитиков с Восточного побережья друзей в среднем больше (табл. 5.2), а если рассматривать только пользователей без ученой степени, то у аналитиков с Восточного побережья друзей снова оказывается в среднем больше!

Таблица 5.2. География и ученая степень пользователей сети

Побережье	Степень	Количество пользователей	Среднее число друзей
Западное	Есть	35	3,1
Восточное	Есть	70	3,2
Западное	Нет	66	10,9
Восточное	Нет	33	13,4

Как только вы начинаете учитывать ученые степени, корреляция движется в обратную сторону! Группировка данных по признаку "восток — запад" скрыла тот факт,

<sup>9</sup> [https://ru.wikipedia.org/wiki/Парадокс\\_Симпсона](https://ru.wikipedia.org/wiki/Парадокс_Симпсона).

<sup>10</sup> Спутывающая переменная (англ. *confounding variable*) — внешняя переменная статистической модели, которая коррелирует (положительно или отрицательно) как с зависимой, так и с независимой переменной. Ее игнорирование приводит к смещению оценки модели. — *Прим. пер.*

что среди аналитиков с Восточного побережья имеется сильная асимметрия в сторону ученых степеней<sup>11</sup>.

Подобный феномен возникает в реальном мире с определенной регулярностью. Главной проблемой является то, что корреляция измеряет связь между двумя переменными *при прочих равных условиях*. Если данным назначать классы случайным образом, что обычно происходит в хорошо поставленном эксперименте, то может оказаться, что допущение "при прочих равных условиях" может иметь место. Но если в основе назначения классов лежит некая многоуровневая схема, то это допущение приводит к неприятным последствиям.

Единственный реальный способ избежать таких неприятностей — это *знать* свои данные и постараться проверить возможные спутывающие факторы. Очевидно, это не всегда можно сделать. Если бы у вас не было данных об уровне образования 200 аналитиков, вы бы просто решили, что есть нечто, присущее людям с Западного побережья, делающее их более общительными.

## Некоторые другие ловушки корреляции

Корреляция, равная нулю, означает отсутствие линейной связи между двумя переменными. Однако могут быть совсем другие виды зависимостей. Например, если:

$$x = [-2, -1, 0, 1, 2]$$

$$y = [2, 1, 0, 1, 2]$$

то переменные  $x$  и  $y$  имеют нулевую корреляцию, но связь между ними определенно существует — каждый элемент  $y$  равен абсолютному значению соответствующего элемента  $x$ . Но в них отсутствует связь, при которой, зная о соотношении переменной  $x_i$  со средним  $\text{mean}(x)$ , можно получить информацию о соотношении переменной  $y_i$  со средним  $\text{mean}(y)$ . Это как раз тот тип связи, который корреляция пытается установить.

Кроме того, корреляция ничего не говорит о величине связи. Переменные:

$$x = [-2, 1, 0, 1, 2]$$

$$y = [99.98, 99.99, 100, 100.01, 100.02]$$

имеют очень хорошую корреляцию, но в зависимости от того, что вы измеряете, вполне возможно, что эта связь не представляет особого интереса.

## Корреляция и причинная зависимость

Наверное, некоторым доводилось слышать утверждение, что "корреляция — это не каузация". Скорее всего, эта фраза принадлежит тому, кто, глядя на данные, которые противоречат какой-то части его мировоззрения, отказывается ставить ее под

<sup>11</sup> *Парадокс Симпсона* или парадокс объединения — статистический эффект, когда при наличии двух групп данных, в каждой из которых наблюдается одинаково направленная зависимость, при объединении этих групп направление зависимости меняется на противоположное. Иллюстрирует неправоту некоторых обобщений (см. [https://ru.wikipedia.org/wiki/Парадокс\\_Симпсона](https://ru.wikipedia.org/wiki/Парадокс_Симпсона)). — *Прим. пер.*

сомнение. Тем не менее, важно понимать, что если  $x$  и  $y$  сильно коррелированы, то это может означать, что либо  $x$  влечет за собой  $y$ , либо  $y$  влечет за собой  $x$ , либо они взаимообусловлены, либо они зависят от какого-то третьего фактора, либо это вовсе ничего не означает.

Рассмотрим связь между списками `num_friends` и `daily_minutes`. Возможно, что чем больше у пользователя друзей на сайте, тем больше он проводит там времени. Это может быть из-за того, что каждый друг ежедневно публикует несколько сообщений, и поэтому чем больше у вас друзей, тем больше вам требуется времени, чтобы оставаться в курсе последних новостей.

Однако возможно и то, что чем больше времени пользователи проводят в спорах на форумах DataSciencester, тем чаще они знакомятся с единомышленниками и становятся друзьями. Иными словами, увеличение времени, проводимого на сайте, *приводит* к росту числа друзей у пользователей.

Третий вариант заключается в том, что пользователи, которые сильнее увлечены наукой о данных, проводят на сайте больше времени (потому что им это интересно) и более активно собирают вокруг себя друзей, специализирующихся в этой области (потому что они не хотят общаться с кем-либо еще).

Один из способов укрепить свою позицию в отношении причинной зависимости заключается в проведении рандомизированных испытаний. Если случайным образом распределить пользователей на две группы (экспериментальную и контрольную) с похожей демографией и предоставить экспериментальной группе несколько иной опыт взаимодействия, то нередко можно получить подтверждение, что разный опыт взаимодействия приводит к разным результатам.

Например, если вы не боитесь, что вас обвинят в экспериментах над вашими пользователями<sup>12</sup>, то вы могли бы случайным образом выделить из числа пользователей сети подгруппу и показывать им только небольшую часть их друзей. Если эта подгруппа в дальнейшем проводила бы на сайте меньше времени, то этот факт вселил бы в вас некоторую уверенность, что наличие большего числа друзей является *причиной* большего количества времени, проводимого на сайте.

## Для дальнейшего изучения

- ◆ Библиотеки SciPy (<http://docs.scipy.org/doc/scipy/reference/stats.html>), pandas (<http://pandas.pydata.org/>) и StatsModels (<http://statsmodels.sourceforge.net/>) поставляются с широким набором статистических функций.
- ◆ Математическая статистика имеет *особое* значение (надо признать, и статистики тоже). Если вы хотите стать хорошим аналитиком данных, то стоит прочесть учебник по математической статистике. Многие из них доступны в Интернете.

---

<sup>12</sup> <http://www.nytimes.com/2014/06/30/technology/facebook-tinkers-with-users-emotions-in-news-feed-experiment-stirring-outcry.html>.

# Теория вероятностей

Законы вероятности, такие истинные в целом и такие ошибочные в частности.

Эдвард Гиббон<sup>1</sup>

Трудно решать задачи в области науки о данных, не имея представления о *теории вероятностей* и ее математическом аппарате. Также как и с обсуждением математической статистики в *главе 5*, здесь мы снова будем много размахивать руками, но обойдем стороной значительную часть технических подробностей.

Для наших целей будем представлять вероятность как способ количественной оценки неопределенности, связанной с *событиями* из некоторого *вероятностного пространства*. Вместо того чтобы вдаваться в технические тонкости по поводу точного определения этих терминов, лучше представьте бросание кубика. Вероятностное пространство состоит из множества всех возможных элементарных исходов (в случае кубика их шесть — по числу граней), а любое подмножество этих исходов представляет собой случайное событие, например, такое как "выпала грань с единицей" или "выпала грань с четным числом".

Вероятность случайного события  $E$  мы будем обозначать  $P(E)$ .

В дальнейшем мы будем применять теорию вероятностей как при построении моделей, так и при их вычислении. Иными словами, она будет использоваться постоянно.

Конечно, при желании можно погрузиться в философию и порассуждать о *сущности* теории вероятностей (этим неплохо заниматься за кружкой пива, например), но мы не будем этого делать.

## Зависимость и независимость

Вообще, два события  $E$  и  $F$  являются *зависимыми*, если какое-то знание о наступлении события  $E$  дает нам информацию о наступлении события  $F$ , и наоборот. В противном случае они *независимые*.

Например, если мы дважды бросаем уравновешенную монету, то знание о том, что в первый раз выпадет орел, не дает никакой информации о том, что орел выпадет и во второй раз. Эти события независимые. С другой стороны, знание о том, что

<sup>1</sup> Эдвард Гиббон (1737–1794) — знаменитый английский историк, автор "Истории упадка и разрушения Римской империи". — *Прим. пер.*

в первый раз выпадет орел, определенно, дает нам информацию о том, выпадут ли решки оба раза. (Если в первый раз выпадет орел, то, конечно же, исключаем случай, когда оба раза выпадают решки.) Оба эти события зависимые.

С точки зрения математики, говорят, что два события  $E$  и  $F$  *независимы*, если вероятность их совместного наступления равна произведению вероятностей их наступления по отдельности:

$$P(E, F) = P(E)P(F).$$

В примере с бросанием монеты вероятность события, что в первый раз выпадет орел, равна  $1/2$ , вероятность события, что оба раза выпадут решки, равна  $1/4$ , а вероятность, что в первый раз выпадет орел и оба раза выпадут решки, равна 0.

## Условная вероятность

Повторим: когда два события  $E$  и  $F$  независимы, то по определению вероятность их совместного наступления равна:

$$P(E, F) = P(E)P(F).$$

Если же они являются зависимыми (и при этом вероятность  $F$  не равна 0), то *условная вероятность* события  $E$  при условии события  $F$  определяется так:

$$P(E|F) = \frac{P(E, F)}{P(F)}.$$

Под этим понимается вероятность наступления события  $E$  при условии, что известно о наступлении события  $F$ .

Эта формула часто приводится к следующему виду:

$$P(E, F) = P(E|F)P(F).$$

В случае, когда  $E$  и  $F$  — независимые события, то можно убедиться, что формула примет вид:

$$P(E|F) = P(E),$$

которая на математическом языке выражает, что знание о наступлении события  $F$  не дает никакой дополнительной информации о наступлении события  $E$ .

Для демонстрации условной вероятности обычно приводят следующий замысловатый пример<sup>2</sup> с семьей, где есть двое детей, чей пол нам не известен. При этом предполагаем, что:

- ◆ каждый ребенок равновероятно является либо мальчиком, либо девочкой;
- ◆ пол второго ребенка не зависит от пола первого.

<sup>2</sup> [https://ru.wikipedia.org/wiki/Парадокс\\_мальчика\\_и\\_девочки](https://ru.wikipedia.org/wiki/Парадокс_мальчика_и_девочки). — Прим. пер.

Тогда вероятность, что оба ребенка не девочки, равна  $1/4$ , вероятность того, что одна девочка и один мальчик, равна  $1/2$ , а вероятность, что обе девочки, равна  $1/4$ .

Для начала узнаем вероятность наступления события, когда оба ребенка — девочки ( $B$ ), при условии, что старший ребенок — девочка ( $G$ ). Используя определение условной вероятности, получим:

$$P(B|G) = \frac{P(B, G)}{P(G)} = \frac{P(B)}{P(G)} = 1/2,$$

поскольку событие  $B$  и  $G$  (оба ребенка — девочки и старший ребенок — девочка) представлено одним событием  $B$ . (Если известно, что оба ребенка — девочки, то верно, что старший ребенок — девочка.)

И, скорее всего, такой результат будет в согласии с нашей интуицией.

А теперь узнаем, какова вероятность события, когда оба ребенка — девочки при условии, что как минимум один ребенок — девочка ( $L$ ). И удивительная вещь — ответ будет отличаться от ранее полученного!

Как и раньше, событие  $B$  и  $L$  (оба ребенка — девочки и не менее одного ребенка — девочка) представлено одним событием  $B$ . Однако получим:

$$P(B|L) = \frac{P(B, L)}{P(L)} = \frac{P(B)}{P(L)} = 1/3.$$

Как такое может быть? Все дело в том, что если известно, что как минимум один ребенок — девочка, то вероятность, что в семье имеется один мальчик и одна девочка, в два раза выше, чем вероятность, что имеются две девочки<sup>3</sup>.

Можно убедиться в правильности рассуждений, "сгенерировав" большое количество семей:

```
# произвольно выбрать мальчика или девочку
```

```
def random_kid():
```

```
    return random.choice(["boy", "girl"])
```

```
# проверка парадокса мальчика и девочки
```

```
both_girls = 0
```

```
older_girl = 0
```

```
either_girl = 0
```

```
random.seed(0)
```

```
for _ in range(10000): # провести эксперимент на совокупности
```

```
    # из 100 000 семей
```

```
    younger = random_kid()
```

<sup>3</sup> Речь идет о парадоксе мальчика и девочки. Когда известно, что в семье как минимум одна девочка, мы автоматически отмечаем вариант с двумя мальчиками. А из того, что оставшиеся три исхода равновероятны, делается вывод, что вероятность "девочка — девочка" равна  $1/3$ . — Прим. пер.

```

older = random_kid()
if older == "girl": # старшая?
    older_girl += 1
if older == "girl" and younger == "girl": # обе?
    both_girls += 1
if older == "girl" or younger == "girl": # любая из двух?
    either_girl += 1

print("P(обе | старшая):", both_girls / older_girl) # 0.514 ~ 1/2
print("P(обе | любая): ", both_girls / either_girl) # 0.342 ~ 1/3

```

## Теорема Байеса

Лучшим другом аналитика данных является *теорема Байеса*<sup>4</sup>, которая позволяет "переставить" условные вероятности местами. Пусть нужно узнать вероятность некоего события  $E$ , зависящего от наступления некоего другого события  $F$ , причем в наличии имеется лишь информация о вероятности события  $F$ , зависящего от наступления события  $E$ . Двукратное применение (в силу симметрии) определения условной вероятности даст формулу Байеса:

$$P(E|F) = \frac{P(E, F)}{P(F)} = \frac{P(F|E)P(E)}{P(F)}$$

Если событие  $F$  разложить на два взаимоисключающих события — событие " $F$  и  $E$ " и событие " $F$  и не  $E$ " — и обозначить "не  $E$ " (т. е.  $E$  не наступает), как  $\neg E$ , тогда:

$$P(F) = P(F, E) + P(F, \neg E),$$

благодаря чему формула приводится к следующему виду:

$$P(E|F) = \frac{P(F|E)P(E)}{P(F|E)P(E) + P(F|\neg E)P(\neg E)}$$

Именно так формулируется теорема Байеса.

Данную теорему часто используют, чтобы продемонстрировать, почему аналитики данных умнее врачей<sup>5</sup>. Представим, что есть некая болезнь, которая поражает 1 из каждых 10 000 человек, и можно пройти обследование, выявляющее эту болезнь, которое в 99% случаев дает правильный результат ("болен", если заболевание имеется, и "не болен" — в противном случае).

Что означает положительный результат обследования? Пусть  $T$  — это событие, что "результат Вашего обследования положительный", а  $D$  — событие, что "у Вас име-

<sup>4</sup> Томас Байес (1702–1761) — английский математик и священник, который первым предложил использование теоремы для корректировки убеждений, основываясь на обновленных данных (см. [https://ru.wikipedia.org/wiki/Теорема\\_Байеса](https://ru.wikipedia.org/wiki/Теорема_Байеса)). — Прим. пер.

<sup>5</sup> См. [https://ru.wikipedia.org/wiki/Теорема\\_Байеса](https://ru.wikipedia.org/wiki/Теорема_Байеса), пример 4. — Прим. пер.



ется заболевание". Тогда согласно теореме Байеса, вероятность наличия заболевания при положительном результате обследования равна:

$$P(D|T) = \frac{P(T|D)P(D)}{P(T|D)P(D) + P(T|\neg D)P(\neg D)}$$

По условию задачи известно, что  $P(T|D) = 0.99$  (вероятность, что заболевший получит положительный результат обследования),  $P(D) = 1/10000 = 0.0001$  (вероятность, что любой человек имеет заболевание),  $P(T|\neg D) = 0.01$  (вероятность, что здоровый человек получит положительный результат обследования) и  $P(\neg D) = 0.9999$  (вероятность, что любое данное лицо не имеет заболевания). Если подставить эти числа в теорему Байеса, то:

$$P(D|T) = 0.98\%$$

т. е. менее 1% людей, которые получают положительный результат обследования, имеют это заболевание на самом деле.



При этом предполагается, что люди проходят обследование более или менее случайно. Если же его проходят только те, у которых имеются определенные симптомы, то вместо этого пришлось бы обуславливать совместным событием "положительный результат обследования и симптомы", в результате чего, возможно, это число оказалось бы намного выше.

В то время как для аналитика данных задача решается простым подсчетом, большинство врачей определяют навскидку, что приблизительно  $P(D|T) = 2$ .

Интуитивно более понятный способ — представить популяцию численностью 1 млн человек. 100 из них ожидаемо имеют заболевание, из которых 99 получили положительный результат обследования. С другой стороны, 999 900 из них ожидаемо не имеют заболевания, из которых 9999 получили положительный результат обследования, вследствие чего можно ожидать, что только 99 из  $(99 + 9999)$  получивших положительный результат обследования имеют это заболевание на самом деле<sup>6</sup>.

## Случайные величины

*Случайная величина* — это переменная, чьим возможным значениям поставлено в соответствие распределение вероятностей. Простая случайная величина равна 1, если подброшенная монета повернется орлом, и 0, если повернется решкой. Более сложная величина может измерять число орлов, наблюдаемых при 10 бросках, или значение, выбираемое из диапазона  $\text{range}(10)$ , где каждое число равновероятное.

<sup>6</sup> Приведенный пример называется *парадоксом теоремы Байеса*, который возникает, если берется редкое явление, такое как туберкулез, например. В этом случае возникает значительная разница в долях больных и здоровых, отсюда и удивительный результат. — *Прим. пер.*

Связанное со случайной величиной распределение вероятностей предоставляет ей вероятности, с которыми она реализует каждое из своих возможных значений. Случайная величина броска монеты равна 0 с вероятностью 0.5 и 1 с той же вероятностью. Случайная величина диапазона  $\text{range}(10)$  имеет распределение, где каждому числу от 0 до 9 поставлена в соответствие вероятность 0.1.

Мы иногда будем говорить о (среднем) ожидаемом значении или *математическом ожидании* случайной величины, которое представляет собой взвешенную сумму произведений каждого ее значения на его вероятность<sup>7</sup>. Среднее ожидаемое значение броска монеты равно  $1/2$  ( $= 0 \cdot 1/2 + 1 \cdot 1/2$ ); среднее ожидаемое значение диапазона  $\text{range}(10)$  равно 4.5.

Случайные величины могут *обуславливаться* событиями точно так же, как и другие события. Пользуясь примером с двумя детьми из *разд. "Условная вероятность"* ранее в данной главе, если  $X$  — это случайная величина, представляющая число девочек, то  $X$  равно 0 с вероятностью  $1/4$ , 1 с вероятностью  $1/2$  и 2 с вероятностью  $1/4$ .

Можно определить новую случайную величину  $Y$ , которая дает число девочек при условии, что как минимум один ребенок — девочка. Тогда  $Y$  равно 1 с вероятностью  $2/3$  и 2 с вероятностью  $1/3$ . А также определить случайную величину  $Z$ , как число девочек при условии, что старший ребенок — девочка, равную 1 с вероятностью  $1/2$  и 2 с вероятностью  $1/2$ .

В рассматриваемых задачах случайные величины будут использоваться по большей части неявным образом, без привлечения к ним особого внимания. Однако если копнуть глубже, то их можно обязательно обнаружить.

## Непрерывные распределения

Бросание монеты соответствует *дискретному распределению*, которое ставит в соответствие положительную вероятность пронумерованным исходам дискретного вероятностного пространства. Однако нередко требуется моделировать распределения на непрерывном пространстве исходов. (Для целей изложения в книге эти исходы всегда будут вещественными числами, хотя в реальной жизни это не всегда так.) Например, *равномерное распределение* назначает *одинаковый вес* всем числам между 0 и 1.

Поскольку между 0 и 1 находится бесконечное количество чисел, то, значит, вес, который оно назначает индивидуальным точкам, должен с неизбежностью быть равен нулю. По этой причине непрерывное распределение расстояния вероятностей представляют плотностью распределения вероятностей (*probability density function*, pdf), также именуемой *дифференциальной функцией распределения* (ДФР), такой, что вероятность наблюдать значение в определенном интервале равна интегралу от дифференциальной функции, взятому в этих пределах.

<sup>7</sup> На практике при анализе выборок математическое ожидание, как правило, неизвестно. Поэтому вместо него используют его оценку — среднее арифметическое. — *Прим. пер.*



Если ваше интегральное исчисление хромает, то более простой способ понять это состоит в том, что если распределение имеет функцию плотности  $f$ , то вероятность наблюдать значение между  $x$  и  $x + h$  приближенно равна  $h \cdot f(x)$  при малых значениях  $h$ .

Функция плотности (или ДФР) равномерного распределения — это всего лишь:

```
# ДФР равномерного распределения
def uniform_pdf(x):
    return 1 if x >= 0 and x < 1 else 0
```

Вероятность, что случайная величина, заданная такой функцией распределения, находится в интервале между 0.2 и 0.3, как и ожидалось, равна 1/10. Функция `random.random()` в Python — это (псевдо)случайная величина с равномерной плотностью.

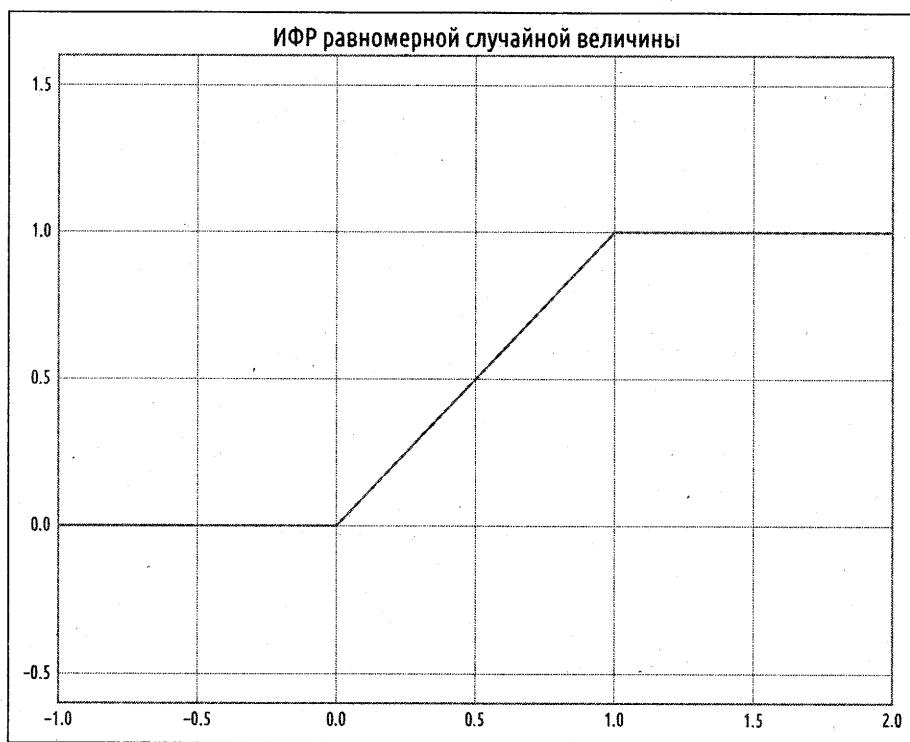


Рис. 6.1. Интегральная функция равномерного распределения

Часто нас больше будет интересовать кумулятивная функция распределения (cumulative distribution function, cdf), также именуемая *интегральной функцией распределения* (ИФР), которая определяет вероятность, что случайная величина меньше или равна некоторому значению. Реализация ИФР равномерного распределения элементарна (рис. 6.1):

```
# ИФР равномерного распределения
def uniform_cdf(x):
```

```

"""возвращает вероятность того, что равномерно распределенная
случайная величина <= x"""
if x < 0: return 0 # величина никогда не бывает меньше 0
elif x < 1: return x # например, P(X <= 0.4) = 0.4
else: return 1 # величина всегда меньше 1

```

## Нормальное распределение

Королем распределений является *нормальное распределение*. Это классическое распределение в форме колокола полностью определяется двумя параметрами:  $\mu$  (мю) — математическим ожиданием (средним значением), и  $\sigma$  (сигмой) — стандартным отклонением. Математическое ожидание указывает на смещение колокола, а стандартное отклонение — на его "ширину" (или масштаб).

Его функция распределения (ДФР) имеет следующий вид:

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right),$$

которую можно реализовать таким образом:

```

# ДФР нормального распределения
def normal_pdf(x, mu=0, sigma=1):
    sqrt_two_pi = math.sqrt(2 * math.pi)
    return (math.exp(-(x-mu) ** 2 / 2 / sigma ** 2) / (sqrt_two_pi * sigma))

```

На рис. 6.2 для наглядности приведены графики некоторых ДФР:

```

xs = [x / 10.0 for x in range(-50, 50)]
plt.plot(xs, [normal_pdf(x, sigma=1) for x in xs], '-', label='мю=0, сигма=1')
plt.plot(xs, [normal_pdf(x, sigma=2) for x in xs], '--', label='мю=0, сигма=2')
plt.plot(xs, [normal_pdf(x, sigma=0.5) for x in xs], ':', label='мю=0, сигма=0.5')
plt.plot(xs, [normal_pdf(x, mu=-1) for x in xs], '-.', label='мю=-1, сигма=1')
plt.legend()
plt.title("Несколько ДФР нормального распределения")
plt.show()

```

При  $\mu=0$  и  $\sigma=1$  оно называется *стандартным нормальным распределением*. Если  $Z$  — это стандартная нормально распределенная случайная величина, то оказывается, что

$$X = \sigma Z + \mu$$

тоже является нормально распределенной, но с математическим ожиданием  $\mu$  и стандартным отклонением  $\sigma$ . И наоборот, если  $X$  — нормально распределенная случайная величина с математическим ожиданием  $\mu$  и стандартным отклонением  $\sigma$ , то

$$Z = \frac{X - \mu}{\sigma}$$

есть стандартная нормально распределенная случайная величина.

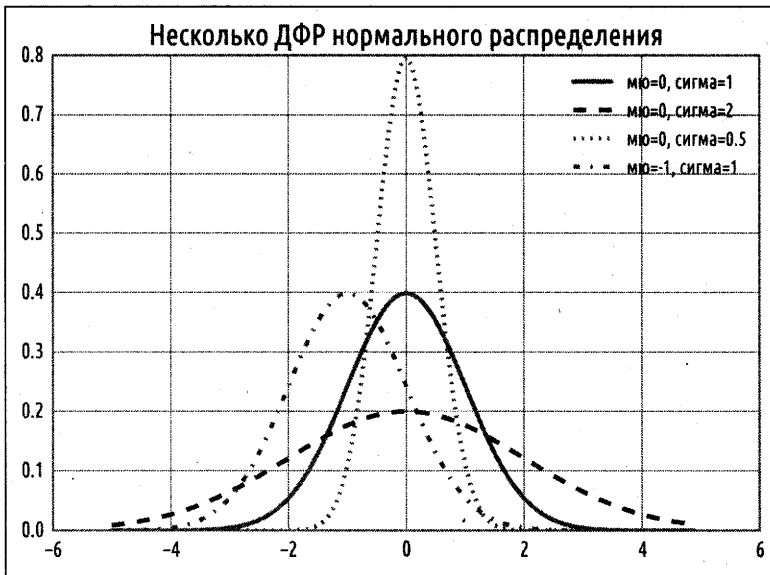


Рис. 6.2. Несколько дифференциальных функций нормального распределения

Интегральную функцию распределения для нормального распределения невозможно написать, пользуясь лишь "элементарными" средствами, однако это можно сделать при помощи питоновской функции интеграла вероятности `math.erf`:

```
# ИФР нормального распределения
def normal_cdf(x, mu=0, sigma=1):
    return (1 + math.erf((x - mu) / math.sqrt(2) / sigma)) / 2
```

И снова построим графики некоторых из них (рис. 6.3):

```
xs = [x / 10.0 for x in range(-50, 50)]
plt.plot(xs, [normal_cdf(x, sigma=1) for x in xs], '-', label='мю=0, сигма=1')
plt.plot(xs, [normal_cdf(x, sigma=2) for x in xs], '--', label='мю=0, сигма=2')
plt.plot(xs, [normal_cdf(x, sigma=0.5) for x in xs], ':', label='мю=0, сигма=0.5')
plt.plot(xs, [normal_cdf(x, mu=-1) for x in xs], '-.', label='мю=-1, сигма=1')
plt.legend(loc=4) # внизу справа
plt.title("Несколько ИФР нормального распределения")
plt.show()
```

Иногда требуется обратить интегральную функцию `normal_cdf`, чтобы найти значение, соответствующее указанной вероятности. Простой способ вычислить обратную функцию отсутствует, однако если учесть, что `normal_cdf` — непрерывная и монотонно возрастающая функция, то можно применить двоичный поиск<sup>8</sup>:

```
# обратная ИФР нормального распределения
# (tolerance - это константа точности)
```

<sup>8</sup> См. [https://ru.wikipedia.org/wiki/Двоичный\\_поиск](https://ru.wikipedia.org/wiki/Двоичный_поиск). — Прим. пер.

```

def inverse_normal_cdf(p, mu=0, sigma=1, tolerance=0.00001):
    """найти приближенную инверсию, используя двоичный поиск"""
    # если не стандартизировано, то стандартизировать и прошкалировать
    if mu != 0 or sigma != 1:
        return mu + sigma * inverse_normal_cdf(p, tolerance=tolerance)

    low_z, low_p = -10.0, 0 # normal_cdf(-10) = (очень близко к) 0
    hi_z, hi_p = 10.0, 1 # normal_cdf(10) = (очень близко к) 1
    while hi_z - low_z > tolerance:
        mid_z = (low_z + hi_z) / 2 # взять середину
        mid_p = normal_cdf(mid_z) # и значение ИФР в этом месте
        if mid_p < p:
            # значение середины все еще слишком низкое, искать выше его
            low_z, low_p = mid_z, mid_p
        elif mid_p > p:
            # значение середины все еще слишком высокое, искать ниже
            hi_z, hi_p = mid_z, mid_p
        else:
            break

    return mid_z

```

Функция многократно делит интервалы пополам, пока не выйдет на точку  $Z$ , которая достаточно близка к требуемой вероятности.

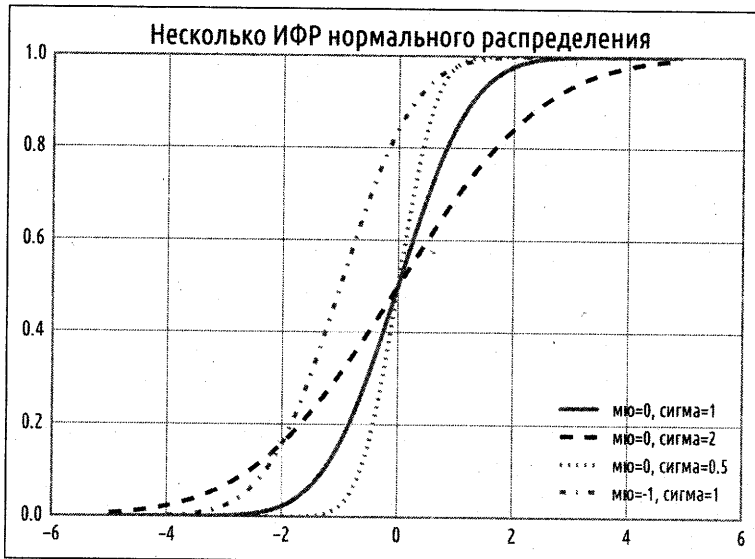


Рис. 6.3. Несколько интегральных функций нормального распределения

## Центральная предельная теорема

Одна из причин распространенности нормального распределения заключается в *центральной предельной теореме* (ЦПТ), согласно которой (по существу) случайная величина, определенная как среднее большого числа независимых и идентично распределенных случайных величин, сама является приближенно нормально распределенной.

В частности, если  $x_1, \dots, x_n$  — случайные величины с математическим ожиданием (средним значением)  $\mu$  и стандартным отклонением  $\sigma$  и если  $n$  большое, то

$$\frac{1}{n}(x_1 + \dots + x_n)$$

— приближенно нормально распределенная величина с математическим ожиданием  $\mu$  и стандартным отклонением  $\sigma/\sqrt{n}$ . Эквивалентным образом (и чаще с большей практической пользой),

$$\frac{(x_1 + \dots + x_n) - \mu n}{\sigma\sqrt{n}}$$

есть приближенно нормально распределенная величина с нулевым математическим ожиданием и стандартным отклонением, равным 1.

Это можно легко проиллюстрировать, обратившись к *биномиально распределенным* случайным величинам, имеющим два параметра —  $n$  и  $p$ . Биномиальная случайная величина  $\text{binomial}(n, p)$  — это просто сумма  $n$  независимых случайных величин с распределением Бернулли  $\text{bernoulli}(p)$ , таких, что значение каждой из них равно 1 с вероятностью  $p$  и 0 с вероятностью  $1 - p$ :

```
# независимое испытание Бернулли, в котором имеется всего
# два случайных исхода (1 и 0) с постоянной вероятностью
def bernoulli_trial(p):
    return 1 if random.random() < p else 0
```

```
# биномиальное распределение
def binomial(n, p):
    return sum(bernoulli_trial(p) for _ in range(n))
```

Математическое ожидание случайной величины с распределением Бернулли  $\text{bernoulli}(p)$  равно  $p$ , ее стандартное отклонение равно  $\sqrt{p(1-p)}$ . Согласно центральной предельной теореме, при больших  $n$  биномиальная случайная величина  $\text{binomial}(n, p)$  приближенно нормально распределена с математическим ожиданием  $\mu = np$  и стандартным отклонением  $\sigma = \sqrt{np(1-p)}$ . На диаграмме легко увидеть сходство обоих распределений:

```
# построить гистограмму
def make_hist(p, n, num_points):
```

```

data = [binomial(n, p) for _ in range(num_points)]

# столбчатая диаграмма, показывающая фактические биномиальные выборки
histogram = Counter(data)
plt.bar([x - 0.4 for x in histogram.keys()],
        [v / num_points for v in histogram.values()],
        0.8,
        color='0.75')

mu = p * n
sigma = math.sqrt(n * p * (1 - p))

# линейный график, показывающий нормальное приближение
xs = range(min(data), max(data) + 1)
ys = [normal_cdf(i + 0.5, mu, sigma) - normal_cdf(i - 0.5, mu, sigma)
      for i in xs]
plt.plot(xs, ys)
plt.title("Биномиальное распределение и его нормальное приближение")
plt.show()

```

Например, при вызове `make_hist(0.75, 100, 10000)` получим диаграмму, как на рис. 6.4.

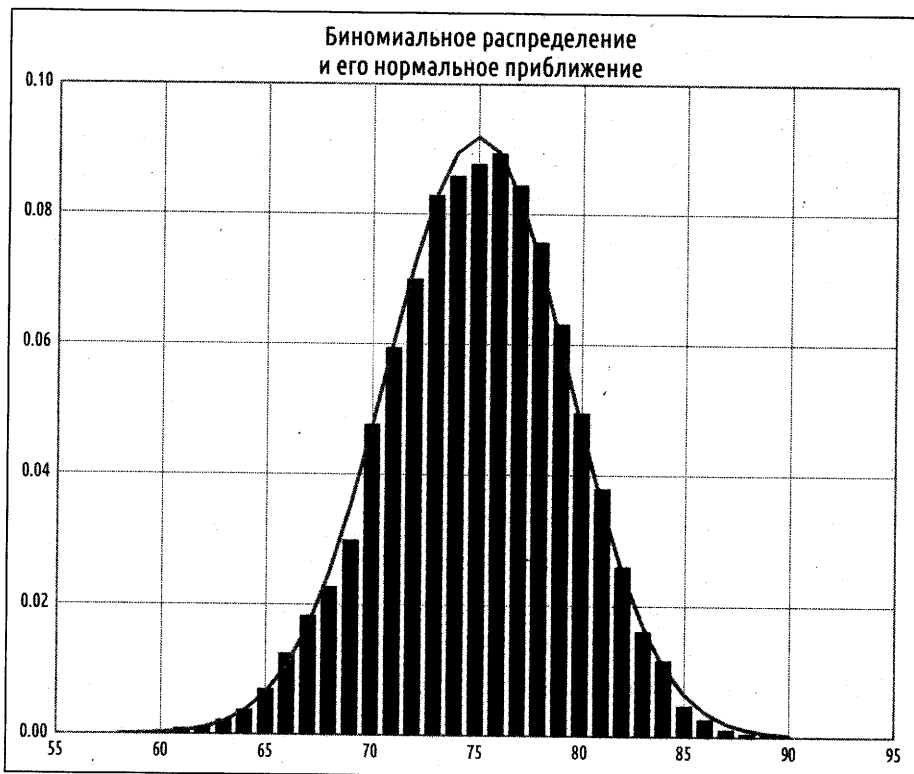


Рис. 6.4. Результат выполнения функции `make_hist`



Мораль этого приближения заключается в том, что если (скажем) нужно узнать вероятность, что уравновешенная монета выпадет орлом более 60 раз из 100 бросков, то это можно вычислить, как вероятность, что `normal(50, 5)` больше 60, и это гораздо легче сделать, чем вычислять ИФР биномиального распределения `binomial(100, 0.5)`. (Хотя в большинстве приложений вы, возможно, воспользуетесь статистическим программным обеспечением, которое с готовностью вычислит любую вероятность, какую пожелаете.)

## Для дальнейшего изучения

- ◆ Библиотека `scipy.stats` (<http://docs.scipy.org/doc/scipy/reference/stats.html>) содержит ДФР и ИФР для большинства популярных распределений вероятностей.
- ◆ В конце главы 5 уже отмечалось, что неплохо было бы изучить учебник по математической статистике. Это же касается и учебника по теории вероятностей.

# Гипотеза и вывод

По-настоящему умного человека характеризует то, что им движет статистика.

Джордж Бернард Шоу<sup>1</sup>

Для чего нужны все эти сведения из математической статистики и теории вероятностей? *Научная* сторона науки о данных часто подразумевает формулировку и проверку *статистических гипотез* о данных и процессах, которые их порождают.

## Проверка статистических гипотез

Аналитикам данных часто требуется выполнить проверку вероятности, что определенная статистическая гипотеза является правильной. Под *статистическими гипотезами* будем понимать утверждения типа "эта монета уравновешена" или "аналитики данных больше предпочитают Python, чем R", или "посетители, скорее всего, покинут страницу, не прочтя содержимого, если разместить на ней всплывающие окна с раздражающей рекламой и крошечной трудноразличимой кнопкой **Заккрыть**". Все эти утверждения могут быть транслированы в *статистики* о данных. В условиях нескольких исходных предположений эти статистики можно рассматривать в качестве наблюдений за случайными величинами с известными распределениями, которые позволяют делать утверждения о возможности, что эти исходные предположения верны.

Классическая трактовка подразумевает наличие главной или *нулевой гипотезы*  $H_0$ , которая представляет некую позицию по умолчанию, и конкурирующей или *альтернативной гипотезы*  $H_1$ , относительно которой мы хотим ее сопоставить. Чтобы принять решение, можно ли отклонить  $H_0$  как ошибочную или принять ее как верную, используют статистики. По-видимому, будет разумнее показать это на примере.

## Пример: бросание монеты

Пусть имеется монета, которую требуется проверить, уравновешена ли она. Для этого делается исходное предположение, что монета имеет некую вероятность  $p$  выпадения орла, и выдвигается нулевая гипотеза о том, что монета уравновешена, т. е.  $p = 0.5$ . Проверим ее, сопоставив с альтернативной гипотезой  $p \neq 0.5$ .

<sup>1</sup> Джордж Бернард Шоу (1856–1950) — ирландский драматург, писатель, романист. — Прим. пер.

В частности, проверка нулевой гипотезы будет состоять в бросании монеты  $n$  раз и подсчете количества орлов  $X$ . Каждый бросок монеты — это испытание Бернулли, где  $X$  — это биномиальная случайная величина  $\text{binomial}(n, p)$ , которую, как мы уже убедились в *главе 6*, можно приближенно выразить при помощи нормального распределения:

```
# аппроксимация биномиальной случайной величины нормальным распределением
def normal_approximation_to_binomial(n, p):
    """находит mu и sigma, которые соответствуют binomial(n, p)"""
    mu = p * n
    sigma = math.sqrt(p * (1 - p) * n)
    return mu, sigma
```

В случае если случайная величина подчиняется нормальному распределению, то для вычисления вероятности, что ее реализованное значение лежит в пределах (или за пределами) определенного интервала, можно воспользоваться функцией `normal_cdf` для вычисления ИФР нормального распределения:

```
# вероятность, что значение нормальной случайной величины лежит
# ниже порогового значения
normal_probability_below = normal_cdf

# оно выше порогового значения, если оно не ниже порогового значения
def normal_probability_above(lo, mu=0, sigma=1):
    return 1 - normal_cdf(lo, mu, sigma)

# оно лежит между, если оно меньше hi, но не ниже lo
def normal_probability_between(lo, hi, mu=0, sigma=1):
    return normal_cdf(hi, mu, sigma) - normal_cdf(lo, mu, sigma)

# оно лежит за пределами, если оно не внутри
def normal_probability_outside(lo, hi, mu=0, sigma=1):
    return 1 - normal_probability_between(lo, hi, mu, sigma)
```

Помимо этого, можно сделать и обратное — найти нехвостовую область или же (симметричный) интервал вокруг среднего значения, на который приходится определенный уровень правдоподобия<sup>2</sup>. Например, если нужно найти интервал с центром в среднем значении, содержащий вероятность 60%, то находим точки отсечения, где верхний и нижний "хвосты" содержат по 20% вероятности каждый (с остатком в 60%):

<sup>2</sup> В статистике значения понятий вероятности и правдоподобия различаются по роли параметра и результата. Вероятность используется для описания результата функции при наличии фиксированного значения параметра. (Если сделать 10 бросков уравновешенной монеты, какова вероятность, что она повернется 10 раз орлом?) Правдоподобие используется для описания параметра функции при наличии ее результата. (Если брошенная 10 раз монета повернулась 10 раз орлом, насколько правдоподобна уравновешенность монеты?) — *Прим. пер.*

```

# верхняя граница нормальной величины
def normal_upper_bound(probability, mu=0, sigma=1):
    """возвращает z, для которого P(Z <= z) = probability"""
    return inverse_normal_cdf(probability, mu, sigma)

# нижняя граница нормальной величины
def normal_lower_bound(probability, mu=0, sigma=1):
    """возвращает z, для которого P(Z >= z) = probability"""
    return inverse_normal_cdf(1 - probability, mu, sigma)

# двусторонние границы нормальной величины
def normal_two_sided_bounds(probability, mu=0, sigma=1):
    """возвращает симметричные (вокруг среднего значения) границы,
    в пределах которых содержится указанная вероятность"""
    tail_probability = (1 - probability) / 2

    # верхняя граница должна иметь значение хвостовой вероятности
    # tail_probability выше ee
    upper_bound = normal_lower_bound(tail_probability, mu, sigma)

    # нижняя граница должна иметь значение хвостовой вероятности
    # tail_probability ниже ee
    lower_bound = normal_upper_bound(tail_probability, mu, sigma)

    return lower_bound, upper_bound

```

В частности, пусть решено сделать  $n=1000$  бросков. Если гипотеза об уравновешенности монеты правильная, то случайная величина  $X$  должна быть распределена приблизительно нормально со средним значением, равным 500, и стандартным отклонением 15.8:

```
mu_0, sigma_0 = normal_approximation_to_binomial(1000, 0.5)
```

На следующем этапе надо принять решение об *уровне значимости*, на котором в дальнейшем и будет сделан вывод о справедливости гипотезы, т. е. насколько мы готовы в результате проверки совершить *ошибку первого рода* ("ложноположительную"), неправильно отклонив  $H_0$ , даже если она правильная. По причинам, затерявшимся в анналах истории, эта готовность часто задана на уровне 5 или 1%. Возьмем за основу 5%.

Рассмотрим проверку, которая отклоняет  $H_0$ , если  $X$  выпадает за пределы границ, заданных следующим образом:

```

# двусторонние границы нормальной случайной величины
normal_two_sided_bounds(0.95, mu_0, sigma_0) # (469, 531)

```

Приняв  $p$  действительно равным 0.5 (т. е.  $H_0$  правильная), имеем всего 5% шансов, что наблюдаемая случайная величина  $X$  будет лежать за пределами этого интерва-

ла, и это в точности соответствует тому уровню значимости, который был нужен. Выражаясь иначе, если  $H_0$  правильная, то приблизительно в 19 случаев из 20 такая проверка будет давать корректный результат.

Нередко также вызывает интерес *мощность* проверки, т. е. вероятность не сделать *ошибку второго рода*, когда не удается отклонить  $H_0$ , даже если она ошибочная. Чтобы ее измерить, следует конкретизировать, что в точности *означает* ошибочность  $H_0$ . (Знание о том, что  $p \neq 0.5$ , не дает обильной информации о распределении  $X$ .) В частности, проверим, что случится, если на самом деле  $p = 0.55$ , благодаря чему свойство уравновешенности монеты чуть смещено в сторону орлов.

В этом случае мощность проверки можно вычислить так:

```
# 95% границы при условии, что p = 0.5
lo, hi = normal_two_sided_bounds(0.95, mu_0, sigma_0)

# фактические mu и sigma при p = 0.55
mu_1, sigma_1 = normal_approximation_to_binomial(1000, 0.55)
```

```
# ошибка второго рода означает: не удалось отклонить нулевую гипотезу;
# это происходит, когда X все еще внутри первоначального интервала
type_2_probability = normal_probability_between(lo, hi, mu_1, sigma_1)
power = 1 - type_2_probability # мощность = 0.887
```

Теперь допустим, что в качестве нулевой гипотезы дано утверждение о том, что монета не смещена в сторону орлов, или  $p \leq 0.5$ . В этом случае нужна *односторонняя* проверка, которая отклоняет нулевую гипотезу, когда  $X$  больше 500, и не отклоняет, когда  $X$  меньше 500. При этом проверка на 5%-м уровне значимости влечет за собой применение функции `normal_probability_below` для нахождения точки отсечения, ниже которой лежит 95% вероятности:

```
# верхняя граница нормальной случайной величины
hi = normal_upper_bound(0.95, mu_0, sigma_0)
# = 526 (< 531, поскольку нужно больше вероятности в верхнем хвосте)

# вероятность ошибки второго рода
type_2_probability = normal_probability_below(hi, mu_1, sigma_1)
power = 1 - type_2_probability # мощность = 0.936
```

Это более мощная проверка, поскольку вместо того, чтобы отвергнуть  $H_0$ , когда  $X$  ниже 469 (что вряд ли вообще произойдет, если альтернативная  $H_1$  правильная), отвергаем нулевую  $H_0$ , когда  $X$  лежит между 526 и 531 (что вполне может произойти, если альтернативная гипотеза  $H_1$  правильная).

## P-значения

Альтернативный подход по сравнению с предыдущим способом проверки предполагает использование *p-значений*<sup>3</sup> или иначе значений уровня *p* значимости. Вместо выбора границ на основе некоторой точки отсечения вероятности вычисляется вероятность (при условии, что  $H_0$  правильная) получения как минимум такого же предельного значения, как и то, которое фактически наблюдалось.

Для двусторонней проверки гипотезы об уравновешенности монеты вычисляем:

```
# двустороннее p-значение
def two_sided_p_value(x, mu=0, sigma=1):
    if x >= mu:
        # если x больше среднего значения, то значения в хвосте больше x
        return 2 * normal_probability_above(x, mu, sigma)
    else:
        # если x меньше среднего значения, то значения в хвосте меньше x
        return 2 * normal_probability_below(x, mu, sigma)
```

При выпадении 530 орлов вызываем:

```
two_sided_p_value(529.5, mu_0, sigma_0) # 0.062
```



Почему 529,5, а не 530? Все дело в поправке на непрерывность<sup>4</sup>. Она отражает тот факт, что вызов функции `normal_probability_between(529.5, 530.5, mu_0, sigma_0)` является более точной оценкой вероятности наблюдать 530 орлов, чем `normal_probability_between(530, 531, mu_0, sigma_0)`.

Поэтому вызов функции `normal_probability_above(529.5, mu_0, sigma_0)` является более точной оценкой вероятности наблюдать как минимум 530 орлов. Отметим, что такой же прием использован в примере, результат выполнения которого показан на рис. 6.4.

Чтобы убедиться, что этот результат является разумной оценкой, проведем модельное испытание:

```
extreme_value_count = 0 # число предельных значений
for _ in range(100000):
    num_heads = sum(1 if random.random() < 0.5 else 0 # подсчитать число орлов
                    for _ in range(1000)) # при 1000 бросках
    if num_heads >= 530 or num_heads <= 470: # подсчитать, как часто
        extreme_value_count += 1 # число 'предельно'

print(extreme_value_count / 100000) # 0.062
```

Поскольку *p*-значение превышает заданный 5%-й уровень значимости, то нулевая гипотеза не отвергается. И напротив, при выпадении 532 орлов *p*-значение будет:

```
two_sided_p_value(531.5, mu_0, sigma_0) # 0.0463
```

<sup>3</sup> Фактически, *p*-значение — это вероятность ошибки при отклонении нулевой гипотезы (ошибки первого рода). — Прим. пер.

<sup>4</sup> См. [https://en.wikipedia.org/wiki/Continuity\\_correction](https://en.wikipedia.org/wiki/Continuity_correction). — Прим. пер.

что меньше 5%-го уровня значимости, и, следовательно, нулевая гипотеза будет отвергнута. Это в точности такая же проверка, что и прежде. Разница лишь в подходе к статистикам.

Верхнее и нижнее  $p$ -значения можно получить аналогичным образом:

```
upper_p_value = normal_probability_above # верхнее p-значение
lower_p_value = normal_probability_below # нижнее p-значение
```

Для проведения односторонней проверки при выпадении 525 орлов получаем:

```
upper_p_value(524.5, mu_0, sigma_0) # 0.061
```

и, следовательно, нулевая гипотеза не будет отвергнута. При выпадении 527 орлов получим:

```
upper_p_value(526.5, mu_0, sigma_0) # 0.047
```

и нулевая гипотеза будет отвергнута.



Перед применением функции `normal_probability_above` для вычисления  $p$ -значений следует обеспечить, чтобы данные были приблизительно нормально распределены. Полно примеров плохих расчетов, когда исходят исключительно из того, что шанс случайного наступления некоторого наблюдаемого события равен одному на миллион, когда на самом деле имеется в виду "шанс при условии, что данные нормально распределены". Если это не так, то расчеты становятся бессмысленными.

Для проверки на соответствие нормальному распределению существуют разнообразные статистические проверки, но даже если просто разместить данные на диаграмме, то это уже станет шагом в правильном направлении.

## Доверительные интервалы

В предыдущем разделе проверялись гипотезы о значении вероятности  $p$  выпадения орлов, которое для "орлов" является *параметром* неизвестного распределения. В этом случае используют третий подход — строят *доверительный интервал* вокруг наблюдаемого значения параметра.

Например, можно оценить вероятность неуравновешенности монеты, обратившись к среднему арифметическому значению случайных величин Бернулли, соответствующих каждому броску монеты — 1 (орел) и 0 (решка). Если наблюдения показывают 525 орлов из 1000 бросков, то оценочно  $p = 0.525$ .

Насколько можно *доверять* такой оценке? Дело в том, что если бы имелось точное значение  $p$ , то согласно центральной предельной теореме (см. разд. "Центральная предельная теорема" главы 6) среднее этих случайных величин с распределением Бернулли должно быть приближенно нормальным с математическим ожиданием  $p$  и стандартным отклонением:

```
math.sqrt(p * (1 - p) / 1000)
```

Но здесь значение  $p$  неизвестно, и поэтому вместо него используется оценка:

```
p_hat = 525 / 1000      # пропорция орлов в выборке (p "с крышкой")
mu = p_hat
sigma = math.sqrt(p_hat * (1 - p_hat) / 1000)  # 0.0158
```

Это не совсем оправданно, и тем не менее, все равно поступают именно так. И используя аппроксимацию нормальным распределением, делаем вывод, что с "уверенностью на 95%" следующий интервал содержит истинный параметр  $p$ :

```
normal_two_sided_bounds(0.95, mu, sigma)      # [0.4940, 0.5560]
```



Это утверждение касается *интервала*, а не  $p$ . Следует понимать его, как утверждение, что, если бы пришлось повторять эксперимент много раз, то в 95% случаев "истинный" параметр (который каждый раз одинаков) будет лежать в пределах наблюдаемого доверительного интервала (который каждый раз может быть разным).

В частности, вывод о том, что монета не уравновешена, не делается, поскольку 0.5 попадает в пределы доверительного интервала.

И напротив, если бы выпало 540 орлов, то было бы следующее:

```
p_hat = 540 / 1000
mu = p_hat
sigma = math.sqrt(p_hat * (1 - p_hat) / 1000)  # 0.0158
normal_two_sided_bounds(0.95, mu, sigma)      # [0.5091, 0.5709]
```

Здесь утверждение, что "монета уравновешена", не лежит в доверительном интервале. (Гипотеза об уравновешенности монеты не проходит проверки, которую, как ожидалось, она должна проходить в 95% случаев, если бы была правильной.)

## Подгонка $p$ -значения

Процедура, которая отвергает нулевую гипотезу только в 5% случаев — по определению, — в 5% случаев будет отвергать нулевую гипотезу ошибочно:

```
# провести эксперимент
def run_experiment():
    """бросить уравновешенную монету 1000 раз,
    True = орлы, False = решки"""
    return [random.random() < 0.5 for _ in range(1000)]

# отвергнуть уравновешенность монеты
def reject_fairness(experiment):
    """используя 5%-е уровни значимости"""
    num_heads = len([flip for flip in experiment if flip])  # число орлов
    return num_heads < 469 or num_heads > 531

random.seed(0)
experiments = [run_experiment() for _ in range(1000)]
```



```
num_rejections = len([experiment      # число отклонений
                       for experiment in experiments
                       if reject_fairness(experiment)])

print(num_rejections) # 46
```

А это означает, что, если задаться целью найти "значимые" результаты, то их обязательно найдешь. Для этого всего лишь нужно сопоставить достаточное количество гипотез с данными, и одна из них почти определенно покажется значимой. Затем, удалив правильные выбросы, в итоге вполне можно получить  $p$ -значение ниже 0.05. (Нечто отдаленно похожее было сделано в разд. "Корреляция" главы 5. Заметили?)

Этот прием, который иногда называют подгонкой  $p$ -значения<sup>5</sup> (или взломом  $p$ -значения,  $p$ -hacking), в некоторой степени является следствием "статистического вывода, заданного рамками  $p$ -значений". По этой теме имеется замечательная статья "Земля круглая"<sup>6</sup>, в которой критикуется такой подход.

Научная объективность, или как говорят "хорошая наука", требует формулировать гипотезы до обращения к данным, очищать данные, не держа в уме гипотезы, и помнить, что  $p$ -значения — это не заменители здравого смысла. (Альтернативный подход рассмотрен в разд. "Байесовский статистический вывод" далее в этой главе.)

## Пример: проведение А/В-тестирования

Одна из ваших первостепенных обязанностей в DataSciencester — заниматься оптимизацией опыта взаимодействия. Под этим эвфемизмом скрываются усилия заставить пользователей щелкать на рекламных объявлениях. Один из рекламных агентов разработал рекламу нового энергетического напитка, предназначенного для аналитиков данных, и директору отдела по рекламе нужна ваша помощь с выбором между двумя рекламными объявлениями: *A* ("вкус отличный!") и *B* ("меньше предвзятости!").

Будучи ученым специалистом, вы решаете провести эксперимент с *A/B-тестированием*<sup>7</sup>, в котором будете случайным образом показывать посетителям сайта одну из двух реклам, при этом отслеживая количество нажатий на каждой из них.

<sup>5</sup> См. <http://www.nature.com/news/scientific-method-statistical-errors-1.14700>.

<sup>6</sup> См. [http://ist-socrates.berkeley.edu/~maccoun/PP279\\_Cohen1.pdf](http://ist-socrates.berkeley.edu/~maccoun/PP279_Cohen1.pdf). Аналогичная статья, в которой рассматривается тема подгонки  $p$ -значения, есть и на русском языке: "Земля круглая ( $p < 0.05$ ): Байес, Фишер и другие" (см. [http://developpsy.ru/publ/alekseev\\_a\\_a\\_zemlja\\_kruglaja\\_p\\_lt\\_0\\_05\\_bajes\\_fisher\\_i\\_drugie/1-1-0-6](http://developpsy.ru/publ/alekseev_a_a_zemlja_kruglaja_p_lt_0_05_bajes_fisher_i_drugie/1-1-0-6)). — Прим. пер.

<sup>7</sup> А/В-тестирование — метод маркетингового исследования, суть которого заключается в том, что контрольная группа элементов сравнивается с набором тестовых групп, в которых один или несколько показателей были изменены, для того чтобы выяснить, какие из изменений улучшают целевой показатель (см. <https://ru.wikipedia.org/wiki/А/В-тестирование>). — Прим. пер.

Если из 1000 потребителей рекламы  $A$  ее выбрали 990 человек, а из 1000 потребителей рекламы  $B$  ее выбрали только 10 человек, то можно быть вполне уверенным, что реклама  $A$  лучше. Но что делать, если разница не такая большая? Как раз здесь и понадобится статистический вывод.

Пусть  $N_A$  людей видят рекламу  $A$  и  $n_A$  из них нажимают на ней. Каждый просмотр рекламы можно представить в виде испытания Бернулли, где  $p_A$  — это вероятность, что кто-то нажмет на рекламе  $A$ . Тогда (если  $N_A$  большое, что так и есть в данном случае) известно, что  $n_A/N_A$  — это приближенно нормальная случайная величина с математическим ожиданием  $p_A$  и стандартным отклонением

$$\sigma_A = \sqrt{p_A(1-p_A)/N_A}.$$

Аналогичным образом,  $n_B/N_B$  — приближенно нормальная случайная величина с математическим ожиданием  $p_B$  и стандартным отклонением  $\sigma_B = \sqrt{p_B(1-p_B)/N_B}$ .

# оценочные параметры

```
def estimated_parameters(N, n):
    p = n / N
    sigma = math.sqrt(p * (1 - p) / N)
    return p, sigma
```

Если допустить, что эти две нормально распределенные случайные величины независимы (что кажется разумным, поскольку отдельные испытания Бернулли обязаны быть такими), то их разность тоже должна быть нормальной с математическим ожиданием  $p_B - p_A$  и стандартным отклонением  $\sqrt{\sigma_A^2 + \sigma_B^2}$ .



Если быть точным, то математика сработает именно так, только если стандартные отклонения *известны*. Здесь же речь идет об их оценках, исходя из выборочных данных, и, следовательно, на самом деле следует использовать  $t$ -распределение. Тем не менее, для достаточно крупных наборов данных их значения настолько близки, что уже нет особой разницы.

Поэтому можно проверить *нулевую гипотезу* о том, что  $p_A$  и  $p_B$  одинаковые (т. е., что  $p_A - p_B = 0$ ), используя следующую статистику:

# статистика  $a/b$ -тестирования

```
def a_b_test_statistic(N_A, n_A, N_B, n_B):
    p_A, sigma_A = estimated_parameters(N_A, n_A)
    p_B, sigma_B = estimated_parameters(N_B, n_B)
    return (p_B - p_A) / math.sqrt(sigma_A ** 2 + sigma_B ** 2)
```

которая должна быть приближенно стандартной нормальной величиной.

Например, если реклама "вкус отличный!" получает 200 откликов из 1000 просмотров, а реклама "меньше предвзятости!" — 180 откликов из такого же количества, то статистика равна:

```
z = a_b_test_statistic(1000, 200, 1000, 180) # -1.14
```

Вероятность наблюдать такую большую разницу при фактически одинаковых средних значениях будет:

```
two_sided_p_value(z) # 0.254
```

что позволит вам сделать вывод о большой разнице. С другой стороны, если реклама "меньше предвзятости!" получит только 150 откликов, то в результате получим:

```
z = a_b_test_statistic(1000, 200, 1000, 150) # -2.94
two_sided_p_value(z) # 0.003
```

и, следовательно, вероятность наблюдать такую большую разницу при одинаково эффективных рекламных объявлениях равна всего 0.003.

## Байесовский статистический вывод

Рассмотренные выше процедуры предполагают выдвижение вероятностных утверждений типа "если нулевая гипотеза правильная, то шанс обнаружить такую-то предельную (экстремальную) статистику равен всего 3%" по поводу *проверок* статистических гипотез.

Альтернативный подход к статистическому выводу связан с рассмотрением в качестве случайных величин самих неизвестных параметров. Аналитик (это вы) начинает с *априорного распределения* для параметров и затем использует наблюдаемые данные и теорему Байеса для получения обновленного *апостериорного распределения* для этих параметров. Вместо вероятностных суждений о проверках делаются вероятностные суждения о самих параметрах.

Например, если неизвестным параметром является вероятность (как в примере с бросанием монеты), то часто априорное распределение вероятности берут из *бета-распределения*<sup>8</sup>, которое размещает всю свою вероятность между 0 и 1:

```
# нормализация
def B(alpha, beta):
    """нормализующая константа, благодаря которой
    сумма вероятностей равна 1"""
    return math.gamma(alpha) * math.gamma(beta) / math.gamma(alpha + beta)

# ДФР для бета-распределенной случайной величины
def beta_pdf(x, alpha, beta):
    if x < 0 or x > 1: # за пределами [0, 1] нет веса
        return 0
    return x ** (alpha - 1) * (1 - x) ** (beta - 1) / B(alpha, beta)
```

<sup>8</sup> Бета-распределение — двухпараметрическое (обычно с произвольными фиксированными параметрами  $\alpha$  и  $\beta$ ) непрерывное распределение. Используется для описания случайных величин, значения которых ограничены конечным интервалом (см. <https://ru.wikipedia.org/wiki/Бета-распределение>). — Прим. пер.

Вообще говоря, это распределение концентрирует свой вес в:

$$\frac{\alpha}{\alpha + \beta}$$

и чем больше значение параметров  $\alpha$  и  $\beta$ , тем "плотнее" распределение.

Например, при  $\alpha=1$  и  $\beta=1$  данное распределение превращается в равномерное (с центром в 0.5 и очень распределенное). При  $\alpha$  намного большем  $\beta$  большая часть веса находится около 1, а при  $\alpha$  намного меньшем  $\beta$  большая часть веса находится около 0. На рис. 7.1 показано несколько различных бета-распределений.

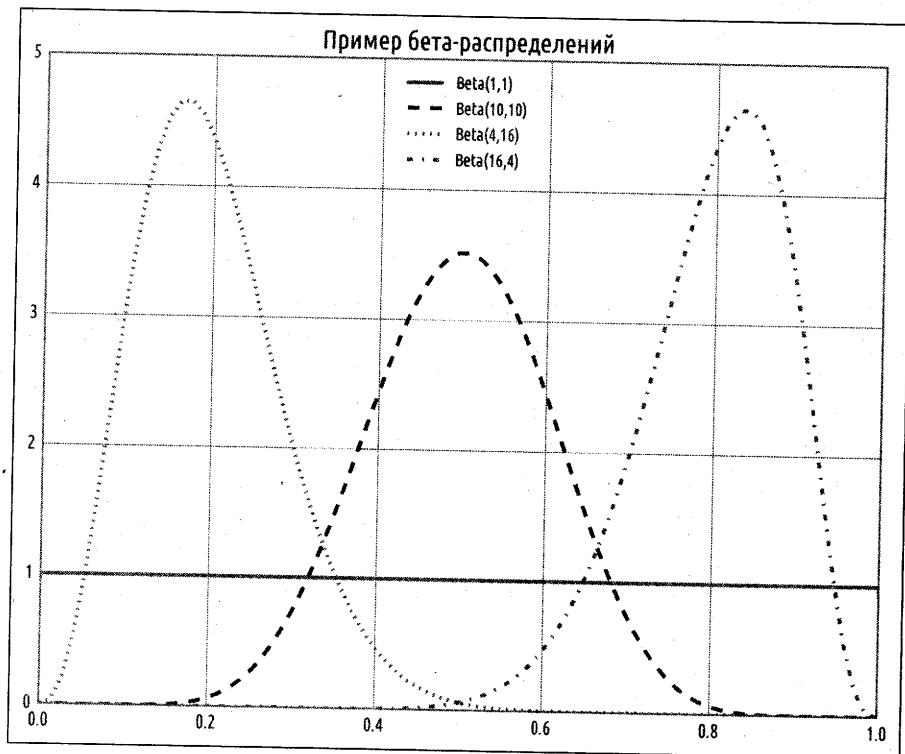


Рис. 7.1. Пример бета-распределений

Пусть задано априорное распределение для  $p$ . Можно не настаивать на уравновешенности монеты и решить, что оба параметра  $\alpha$  и  $\beta$  равны 1. Или же быть абсолютно убежденным в том, что монета будет выпадать орлом в 55% случаев, и потому выбрать  $\alpha=55$ , а  $\beta=45$ .

Затем монету бросают  $n$  раз и подсчитывают  $h$  орлов и  $t$  решек. Согласно теореме Байеса (и некоторым математикам, которых было бы утомительно тут перечислять), апостериорное распределение для  $p$  снова будет бета-распределением, но с параметрами  $\alpha + h$  и  $\beta + t$ .



Неслучайно апостериорное распределение снова подчинено бета-распределению. Число орлов задано биномиальным распределением, а бета-распределение — это априорное распределение, сопряженное с биномиальным<sup>9</sup>. И поэтому, когда априорное бета-распределение обновляется, используя наблюдения из соответствующего биномиального, то в итоге получается апостериорное бета-распределение.

Пусть монета брошена 10 раз, и выпало только 3 орла.

Если начать с равномерного априорного распределения (в некотором смысле, отказываясь отстаивать уравновешенность монеты), то апостериорное распределение будет  $\text{beta}(4, 8)$ , с центром вокруг 0.33. Поскольку все вероятности рассматриваются равновероятными, то предварительная догадка будет находиться где-то поблизости от наблюдаемой вероятности.

Если же начать с  $\text{beta}(20, 20)$  (выражая уверенность, что монета примерно уравновешена), то апостериорное распределение будет  $\text{beta}(23, 27)$  с центром вокруг 0.46, свидетельствуя о пересмотре степени уверенности в сторону того, что, возможно, результаты подбрасывания монеты слегка смещены в сторону решек.

А если начать с  $\text{beta}(30, 10)$  (выражая уверенность, что результаты подбрасывания монеты смещены в сторону выпадения орлов в 75% случаев), то апостериорное распределение будет  $\text{beta}(33, 17)$  с центром вокруг 0.66. В этом случае степень уверенности в смещенности в сторону орлов все еще присутствует, но уже не такая сильная, как первоначально. Эти три апостериорных распределения показаны на рис. 7.2.

Если продолжить бросать монету, то априорное распределение будет значить все меньше и меньше, пока в итоге не получится (почти) одинаковое апостериорное распределение, неважно, с какого априорного распределения оно начиналось.

Например, не имеет значения, насколько, по первоначальной мысли, монета была смещена, поскольку будет трудно поддерживать эту уверенность после выпадения 1000 орлов из 2000 бросков (если, конечно, не быть экстремалом, который решает выбрать, скажем,  $\text{beta}(1000000, 1)$  в качестве априорного бета-распределения).

Интересно, что на основе байесовского вывода можно делать вероятностные утверждения о гипотезах, типа: "исходя из априорного распределения и наблюдавшихся данных, выпадение орлов с вероятностью между 49 и 51% имеет правдоподобие, равное всего 5%". С философской точки зрения такое утверждение сильно отличается от утверждений, типа "если монета уравновешена, то можно ожидать появления таких-то предельных значений только в 5% случаев".

Применение байесовского вывода для проверки статистических гипотез считается несколько противоречивым решением частично потому, что его математический аппарат может становиться весьма запутанным, и частично из-за субъективной природы выбора априорного распределения. Мы больше не будем пользоваться им в этой книге, но знать о нем стоит.

<sup>9</sup> См. [http://www.johndcook.com/blog/conjugate\\_prior\\_diagram/](http://www.johndcook.com/blog/conjugate_prior_diagram/).

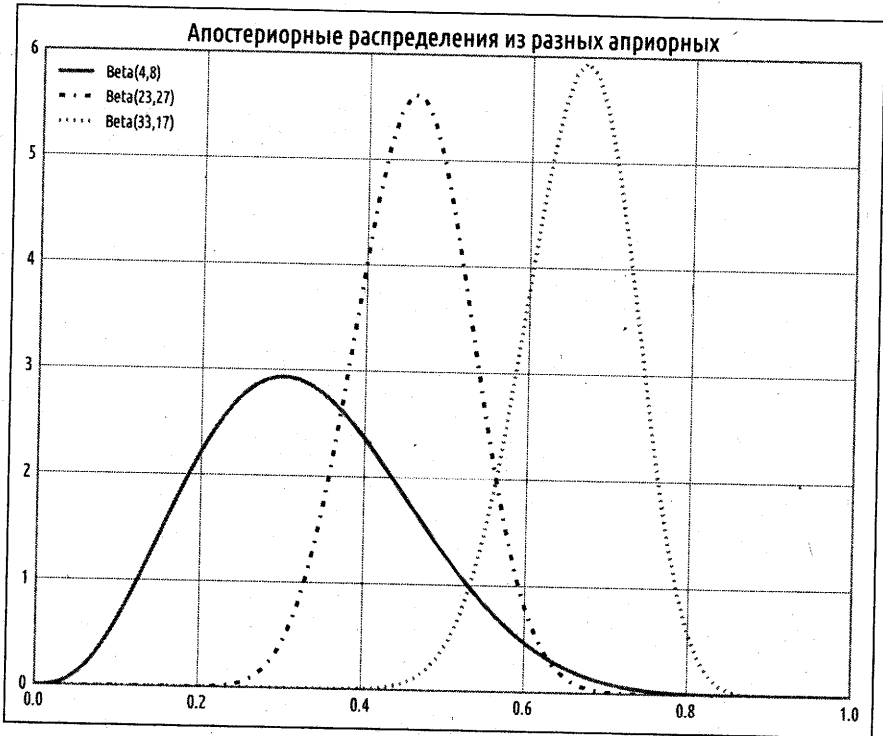


Рис. 7.2. Апостериорные распределения, полученные из разных априорных

## Для дальнейшего изучения

- ◆ В данной главе была затронута лишь незначительная часть того, что необходимо знать о статистическом выводе. В книгах, рекомендованных для прочтения в конце главы 5, данная тема обсуждается гораздо глубже.
- ◆ Образовательная онлайн-платформа Coursera<sup>10</sup> предлагает курс по анализу данных и статистическому выводу (<https://www.coursera.org/specializations/statistics>), который охватывает многие из упомянутых тем.

<sup>10</sup> Coursera — проект в сфере массового онлайн-образования, основанный профессорами информатики Стэнфордского университета, США, предлагающий бесплатные онлайн-курсы от специалистов со всего мира на разных языках. — Прим. пер.

# Градиентный спуск

Те, кто бахвалятся своим происхождением, хвастаются своим долгом перед другими.

Сенека<sup>1</sup>

Решая задачи в области науки о данных, мы будем пытаться подобрать наилучшую модель, соответствующую конкретной ситуации. При этом обычно под "наилучшей" подразумевается модель, которая "минимизирует ошибку модели" либо "максимизирует правдоподобие данных". Другими словами, она должна обеспечивать решение некоторой оптимизационной задачи.

Следовательно, необходимо уметь решать ряд оптимизационных задач. И в данном случае придется их решать с чистого листа. Наш подход будет основываться на *методе градиентного спуска* (gradient descent), который вполне укладывается в трактовку обсуждаемых тем с чистого листа. Рассмотрение данной темы само по себе не является каким-то особо захватывающим занятием, однако оно даст возможность заниматься увлекательными вещами на протяжении всей книги. Так что держитесь рядом.

## Идея в основе метода градиентного спуска

Пусть имеется некая функция  $f$ , которая в качестве входящего аргумента принимает вектор из вещественных чисел и возвращает единственное вещественное число. Вот одна из таких простых функций:

```
# пример целевой функции: сумма квадратов
def sum_of_squares(v):
    """вычисляет сумму квадратов элементов вектора v"""
    return sum(v_i ** 2 for v_i in v)
```

Нередко требуется найти максимум (или минимум) такой функции. Другими словами, надо найти такой входящий аргумент  $v$ , который на выходе дает наибольшее (или наименьшее) возможное значение.

Для таких функций *градиент* (в дифференциальном исчислении это вектор, составленный из частных производных<sup>2</sup>) задает для входящего аргумента направление

<sup>1</sup> Луций Анней Сенека (IV в. до н. э.) — римский философ-стоик, поэт и государственный деятель. Здесь обыгрывается слово *descent*, которое можно понять, как происхождение и как спуск. — *Прим. пер.*

<sup>2</sup> *Частная производная* — предел отношения приращения функции по выбранной переменной к приращению этой переменной при стремлении этого приращения к нулю (см. [https://ru.wikipedia.org/wiki/Частная\\_производная](https://ru.wikipedia.org/wiki/Частная_производная)). — *Прим. пер.*

наибыстрейшего возрастания функции (если сомневаетесь, что еще помните дифференциальное исчисление, то лучше просто поверить на слово либо заглянуть в Интернет<sup>3</sup>).

Поэтому один из подходов к максимизации функции состоит в том, чтобы выбрать произвольную отправную точку, вычислить градиент, сделать малый шаг в направлении градиента (т. е. в направлении, которое заставляет функцию расти быстрее всего) и повторить итерацию с новой отправной точки. Аналогичным образом можно попытаться минимизировать функцию, делая малые шаги в *противоположном* направлении (или в направлении антиградиента), как показано на рис. 8.1.

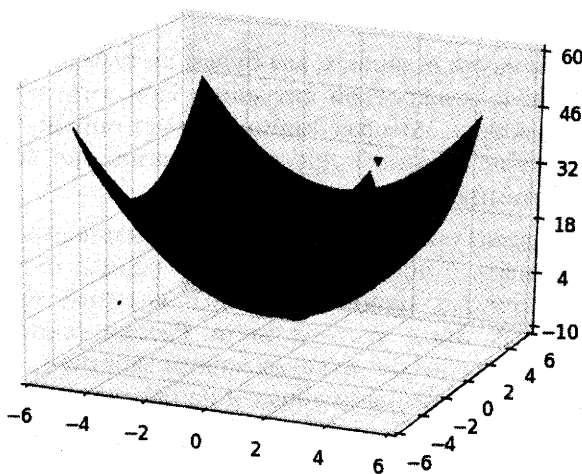


Рис. 8.1. Поиск минимума методом градиентного спуска



Если функция имеет уникальный глобальный минимум, то этот алгоритм, вероятно, его найдет. Если же (локальных) минимумов несколько, то он может "найти" какой-то неправильный из них. В этом случае можно повторить процедуру с нескольких отправных точек. При отсутствии минимума процедура вполне может войти в бесконечный цикл.

## Вычисление градиента

Если  $f$  — это функция одной переменной, то ее производная в точке  $x$  измеряет скорость изменения  $f(x)$  при очень малом изменении в  $x$ . Определяется как предел отношения приращений:

```
# отношение приращений
def difference_quotient(f, x, h):
    return (f(x + h) - f(x)) / h
```

при стремлении  $h$  к нулю.

(Многие из тех, кто изучает исчисления, озадачены таким определением предела. Здесь мы слухавим и просто скажем, что он означает именно то, что вы думаете.)

<sup>3</sup> <https://ru.wikipedia.org/wiki/Градиент>. — Прим. пер.



Производная — это наклон касательной в  $(x, f(x))$ , в то время как отношение приращения — это наклон "не совсем" касательной, которая проходит через  $(x + h, f(x + h))$ . По мере уменьшения  $h$  "не совсем" касательная становится все ближе и ближе к касательной (рис. 8.2).

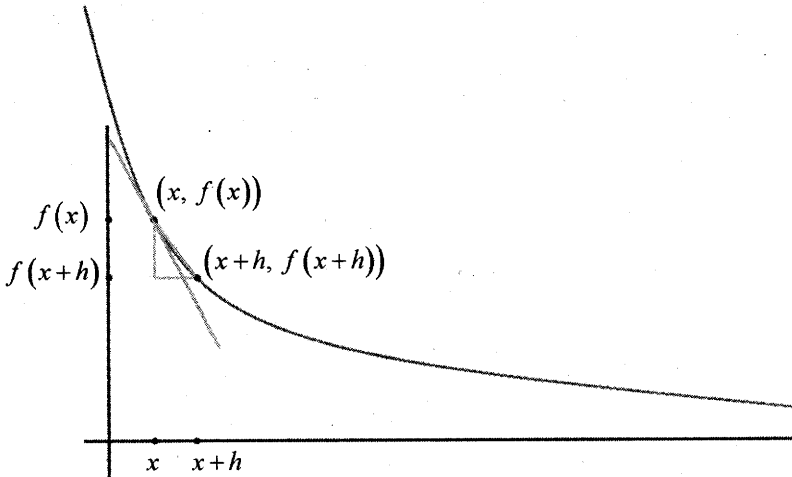


Рис. 8.2. Аппроксимация производной при помощи отношения приращений

Для многих функций достаточно легко выполнить точное вычисление производных. Например, функция возведения в степень `square`:

```
def square(x):
    return x * x
```

имеет производную:

```
def derivative(x):
    return 2 * x
```

в чем, если есть желание, можно убедиться, вычислив отношение приращений явным образом и взяв предел.

А если не получилось (или нет желания) найти градиент? Несмотря на то, что в Python невозможно непосредственно брать пределы, можно выполнить оценку производных путем вычисления отношения приращений для сколь угодно малого положительного числа  $\epsilon$ . На рис. 8.3 показаны результаты одного такого оценивания<sup>4</sup>:

<sup>4</sup> В Python 3 (3.5.2) функция `map` в методе `plot()` может вызывать исключение, поэтому вместо нее следует использовать преобразование последовательности в список (`list(map(...))`) либо генератор последовательности в виде списка, например,

```
plt.plot(xs, [2*x for x in xs], 'rx', label='факт')
plt.plot(xs, [derivative_estimate(x) for x in xs], 'b+', label='Оценка')
```

— Прим. пер.

```

# оценка производной
derivative_estimate = partial(difference_quotient, square, h=0.00001)

# построить диаграмму, чтобы показать, что фактические производные
# и их приближения в сущности одинаковые
import matplotlib.pyplot as plt
xs = range(-10,10)
plt.title("Фактические производные и их оценки в сравнении")
plt.plot(xs, map(derivative, xs), 'rx', label='Факт') # красный x
plt.plot(xs, map(derivative_estimate, xs), 'b+', label='Оценка') # синий +
plt.legend(loc=9)
plt.show()

```

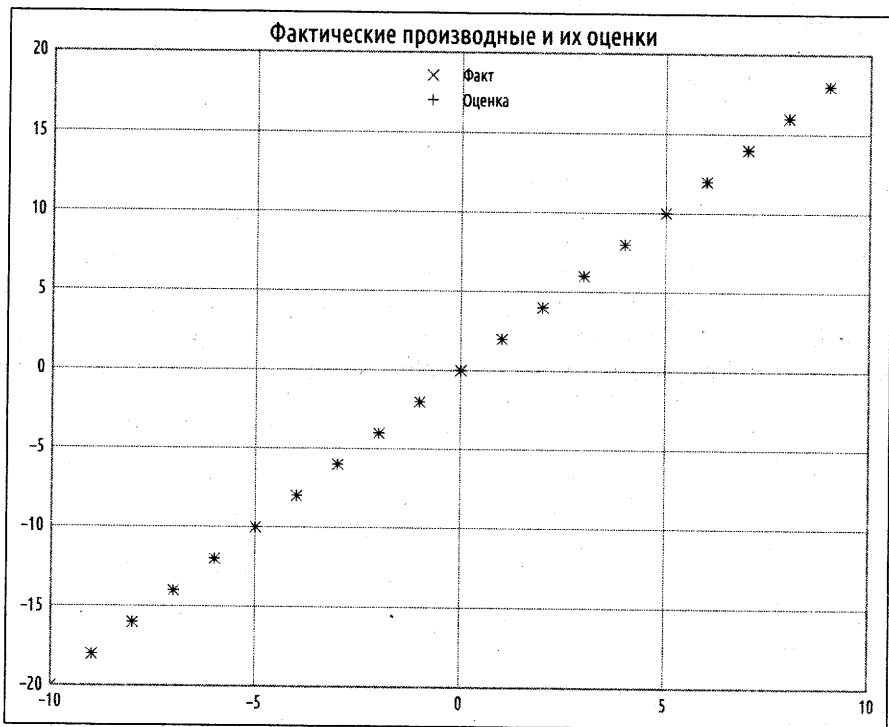


Рис. 8.3. Качество аппроксимации отношением приращения

Если  $f$  — это функция многих переменных, то она имеет кратные *частные производные*, каждая из которых показывает скорость изменения  $f$  при незначительном изменении лишь в одной из входящих переменных.

Ее  $i$ -я частная производная вычисляется как функция одной  $i$ -й переменной при прочих фиксированных переменных:

```

# частное отношение приращений
def partial_difference_quotient(f, v, i, h):

```

```

"""вычислить i-е частное отношение приращений функции f в векторе v"""
w = [v_j + (h if j == i else 0) # прибавить h только к i-му элементу v
      for j, v_j in enumerate(v)]

return (f(w) - f(v)) / h

```

После чего можно таким же образом вычислить градиент:

```

# оценить градиент
def estimate_gradient(f, v, h=0.00001):
    return [partial_difference_quotient(f, v, i, h)
            for i, _ in enumerate(v)]

```



Основной недостаток оценки при помощи отношения приращений заключается в том, что она является вычислительно ресурсоемкой. Если  $v$  имеет длину  $n$ , то функция `estimate_gradient` должна вычислить  $f$  для  $2n$  различных входящих значений. В случае многократной оценки градиентов будет выполняться большой объем избыточной работы.

## Использование градиента

Легко убедиться, что функция суммы квадратов `sum_of_squares` минимальна, когда ее входящий вектор  $v$  состоит из нулей. Претворимся, что мы не знаем об этом. Будем использовать градиенты для нахождения минимума среди всех трехмерных векторов. Выберем произвольную отправную точку и маленькими шагами начнем двигаться в направлении антиградиента, пока не достигнем точки, где градиент очень мал:

```

# сделать шаг градиента
def step(v, direction, step_size):
    """двигаться с шаговым размером step_size в направлении от v"""
    return [v_i + step_size * direction_i
            for v_i, direction_i in zip(v, direction)]

# градиент суммы квадратов
def sum_of_squares_gradient(v):
    return [2 * v_i for v_i in v]

# выбрать произвольную отправную точку
v = [random.randint(-10,10) for i in range(3)]

tolerance = 0.0000001 # константа точности расчета

while True:
    gradient = sum_of_squares_gradient(v) # вычислить градиент в v
    next_v = step(v, gradient, -0.01) # сделать отрицательный шаг градиента
    if distance(next_v, v) < tolerance: # остановиться при достижении
        break # приемлемого уровня (т. е. схождения)
    v = next_v # продолжить, если нет

```

Выполнив этот скрипт, можно убедиться, что он всегда завершается с вектором  $v$ , очень близким к  $[0, 0, 0]$ . При этом, чем меньше задана константа точности расчета `tolerance`, тем ближе он будет к нулевому вектору.

## Выбор оптимального размера шага

Если причина движения против градиента понятна, то насколько далеко двигаться — не совсем. На самом деле, выбор оптимального размера шага больше касается практического, чем научного аспекта науки о данных. Популярные варианты следующие:

- ◆ применение постоянного размера шага;
- ◆ постепенное дробление шага во времени;
- ◆ на каждом шаге выбор размера шага, который минимизирует значение целевой функции (метод наискорейшего спуска).

Последний вариант выглядит наиболее предпочтительным, но на практике тоже требует больших объемов вычислений. Его можно приближенно выразить, пробуя разные размеры шага и выбирая тот, который приводит к наименьшему значению целевой функции:

```
# размеры шага
step_sizes = [100, 10, 1, 0.1, 0.01, 0.001, 0.0001, 0.00001]
```

Некоторые размеры шага могут привести к неправильным входящим аргументам для целевой функции, поэтому необходимо создать функцию-обертку с безопасным применением аргументов, которая для неправильных входящих аргументов возвращает бесконечность (не являющуюся минимумом ни для чего):

```
# безопасная версия
def safe(f):
    """вернуть новую функцию, одинаковую с f, за исключением того,
    что она возвращает бесконечность всякий раз, когда f выдает ошибку"""
    def safe_f(*args, **kwargs):
        try:
            return f(*args, **kwargs)
        except:
            return float('inf') # в Python так обозначается бесконечность
    return safe_f
```

## Собираем все вместе

В общем случае имеется некая целевая функция `target_fn`, которую требуется минимизировать, а также ее градиент `gradient_fn`<sup>5</sup>. К примеру, такой градиент `gradient_fn` может представлять ошибки в модели, как функцию своих параметров,

<sup>5</sup> Градиент целевой функции — это вектор, характеризующий направление и скорость изменения целевой функции. Он определяется ее частными производными по каждой переменной. — Прим. пер.

и тогда мы могли бы отыскать параметры, которые минимизируют ошибки настолько, насколько это возможно.

Далее, пусть каким-то образом выбрано исходное значение для вектора параметров `theta_0`. Тогда метод градиентного спуска реализуется следующим образом:

```
# пакетная минимизация
def minimize_batch(target_fn, gradient_fn, theta_0, tolerance=0.000001):
    """использует градиентный спуск для нахождения вектора theta, который
        минимизирует целевую функцию target_fn"""
    step_sizes = [100, 10, 1, 0.1, 0.01, 0.001, 0.0001, 0.00001]

    theta = theta_0                # установить theta в начальное значение
    target_fn = safe(target_fn)    # безопасная версия
                                   # целевой функции target_fn
    value = target_fn(theta)       # минимизируемое значение

    while True:
        gradient = gradient_fn(theta)
        next_thetas = [step(theta, gradient, -step_size)
                        for step_size in step_sizes]

        # выбрать то, которое минимизирует функцию ошибок
        next_theta = min(next_thetas, key=target_fn)
        next_value = target_fn(next_theta)

        # остановиться, если функция сходится (стремится к пределу)
        if abs(value - next_value) < tolerance: # меньше константы точности?
            return theta
        else:
            theta, value = next_theta, next_value
```

Эта функция названа пакетной минимизацией (`minimize_batch`), потому что на каждом шаге градиента, или итерации, она пересматривает весь набор данных целиком (поскольку целевая функция `target_fn` возвращает ошибку на всем наборе данных). В следующем разделе мы рассмотрим альтернативный подход, который на каждом шаге рассматривает только одну точку данных.

Иногда, напротив, надо *максимизировать* функцию. Это можно сделать путем минимизации ее отрицания (которое имеет соответствующий отрицательный градиент):

```
# отрицание результата на выходе
def negate(f):
    """вернуть функцию, которая для любого входящего x возвращает -f(x)"""
    return lambda *args, **kwargs: -f(*args, **kwargs)

# отрицание списка результатов на выходе
def negate_all(f):
```

```

"""То же самое, когда f возвращает список чисел"""
return lambda *args, **kwargs: [-y for y in f(*args, **kwargs)]

# пакетная максимизация путем минимизации отрицания
def maximize_batch(target_fn, gradient_fn, theta_0, tolerance=0.000001):
    return minimize_batch(negate(target_fn),
                          negate_all(gradient_fn),
                          theta_0,
                          tolerance)

```

## Стохастический градиентный спуск

Как уже упоминалось ранее, метод градиентного спуска, как правило, будет применяться для выбора параметров модели таким образом, чтобы минимизировать некое представление об ошибке<sup>6</sup>. Подход на основе предыдущей пакетной версии метода требует на каждом шаге градиента делать прогноз и вычислять градиент на всем наборе данных, в результате чего на вычисление каждого шага уходит много времени.

Функции ошибок обычно имеют свойство аддитивности, т. е. прогнозная ошибка на всем наборе данных является суммой прогнозных ошибок для каждой точки данных.

В этом случае можно применить *метод стохастического градиентного спуска* (stochastic gradient descent), который за одну итерацию цикла вычисляет градиент (и делает шаг) только на одной точке. Он многократно просматривает данные, пока не достигнет точки останова.

Во время каждого цикла просмотр данных выполняется в случайном порядке:

```

# перемешать индексы
def in_random_order(data):
    """Генератор, который возвращает элементы данных в случайном порядке"""
    indexes = [i for i, _ in enumerate(data)] # создать список индексов
    random.shuffle(indexes) # перемешать данные и
    for i in indexes: # вернуть в этом порядке
        yield data[i]

```

И шаг градиента делается для каждой точки данных. Этот подход не исключает возможности попадания в бесконечный цикл, находясь в окрестностях минимума, поэтому, как только улучшение результата прекращается, размер шага уменьшается, и в конечном счете алгоритм завершается:

<sup>6</sup> Функции `minimize_batch` и `maximize_stochastic` используются в разд. "Снижение размерности" главы 10 и "Применение модели" главы 16. Функция `minimize_stochastic` используется в разд. "Применение метода градиентного спуска" главы 14 и в разд. "Подбор модели" и "Регуляризация" главы 15. — Прим. пер.

*# стохастическая минимизация*

```
def minimize_stochastic(target_fn, gradient_fn, x, y, theta_0, alpha_0=0.01):
```

```
    data = zip(x, y)
```

```
    theta = theta_0 # первоначальная гипотеза
```

```
    alpha = alpha_0 # первоначальный размер шага
```

```
    min_theta, min_value = None, float("inf") # минимум на этот момент
```

```
    iterations_with_no_improvement = 0
```

```
# остановиться, если достигли 100 итераций без улучшений
```

```
while iterations_with_no_improvement < 100:
```

```
    value = sum(target_fn(x_i, y_i, theta) for x_i, y_i in data)
```

```
    if value < min_value:
```

```
        # если найден новый минимум, то запомнить его
```

```
        # и вернуться к первоначальному размеру шага
```

```
        min_theta, min_value = theta, value
```

```
        iterations_with_no_improvement = 0
```

```
        alpha = alpha_0
```

```
    else:
```

```
        # в противном случае улучшений нет,
```

```
        # поэтому пытаемся сжать размер шага
```

```
        iterations_with_no_improvement += 1
```

```
        alpha *= 0.9
```

```
# и делаем шаг градиента для каждой из точек данных
```

```
for x_i, y_i in in_random_order(data):
```

```
    gradient_i = gradient_fn(x_i, y_i, theta)
```

```
    theta = vector_subtract(theta,
```

```
                            scalar_multiply(alpha, gradient_i))
```

```
return min_theta
```

Стохастическая версия обычно работает намного быстрее, чем пакетная. Разумеется, максимизирующая версия тоже понадобится:

*# стохастическая максимизация*

```
def maximize_stochastic(target_fn, gradient_fn, x, y, theta_0, alpha_0=0.01):
```

```
    return minimize_stochastic(negate(target_fn),
```

```
                               negate_all(gradient_fn),
```

```
                               x, y, theta_0, alpha_0)
```

## Для дальнейшего изучения

- ◆ Продолжайте читать! Мы будем пользоваться градиентным спуском для решения задач из остальной части книги.
- ◆ На данный момент некоторые, несомненно, уже почувствовали усталость от моих советов по поводу чтения учебников. Чтобы как-то утешить, вот книга

"Active Calculus" ("Активный математический анализ", <http://scholarworks.gvsu.edu/books/10/>), которая, по-видимому, будет понятнее учебников по математическому анализу, на которых учился я.

- ◆ Библиотека scikit-learn располагает модулем стохастического градиентного спуска (<http://scikit-learn.org/stable/modules/sgd.html>), который по некоторым позициям не такой общий, как представленный в книге, и более общий по другим параметрам. На самом деле в большинстве реальных ситуаций вы, наверняка, будете пользоваться библиотеками, в которых об оптимизации уже позаботились заранее, и вам не придется делать это самим (кроме тех случаев, когда она работает некорректно, что в один прекрасный день неизбежно случится).



# Сбор данных

Писал я ее ровно три месяца, обдумывал ее содержание три минуты, а материал для нее собирал всю жизнь.

Ф. Скотт Фицджеральд<sup>1</sup>

Чтобы быть аналитиком данных, прежде всего нужны сами данные. Специалисты, работающие в этой области, тратят ошеломляюще огромную часть времени на сбор, очистку и преобразование данных. Естественно, в случае крайней необходимости всегда есть возможность внести данные вручную (или поручить это своим миньонам, если они есть), но обычно такая работа не самое лучшее применение собственному времени. В этой главе мы обратимся к различным способам и форматам передачи данных в Python.

## Объекты *stdin* и *stdout*

Если выполнять сценарии Python из командной строки, то можно *передавать* через них данные, пользуясь для этого объектами `sys.stdin` и `sys.stdout`. Например, вот сценарий, который считывает строки текста и выдает обратно те, которые соответствуют регулярному выражению:

```
# egrep.py
import sys, re

# sys.argv - список аргументов командной строки
# sys.argv[0] - имя самой программы
# sys.argv[1] - регулярное выражение, указываемое в командной строке
regex = sys.argv[1]

# для каждой строки, переданной сценарию
for line in sys.stdin:
    # если она соответствует регулярному выражению regex,
    # записать ее в stdout
    if re.search(regex, line): sys.stdout.write(line)
```

А этот сценарий подсчитывает количество полученных строк и возвращает итоговый результат:

---

<sup>1</sup> Фрэнсис Скотт Фицджеральд (1896–1940) — американский писатель, крупнейший представитель так называемого "потерянного поколения" в литературе (см. [https://ru.wikipedia.org/wiki/Фитцджеральд,\\_Фрэнсис\\_Скотт](https://ru.wikipedia.org/wiki/Фитцджеральд,_Фрэнсис_Скотт)). — Прим. пер.

```
# подсчет строк (line_count.py)
import sys

count = 0
for line in sys.stdin:
    count += 1

# результат выводится на консоль sys.stdout
print(count)
```

Этот сценарий в дальнейшем можно использовать для подсчета строк файла, в которых содержатся числа. В Windows это выглядит следующим образом:

```
type SomeFile.txt | python egrep.py "[0-9]" | python line_count.py
```

В UNIX-подобной системе это выглядит несколько иначе:

```
cat SomeFile.txt | python egrep.py "[0-9]" | python line_count.py
```

Символ конвейерной передачи "|" означает, что надо "использовать выход команды слева, как вход для команды справа". Таким способом можно создавать достаточно подробные конвейеры обработки данных.



В Windows эту команду можно использовать без ссылки на Python:

```
type SomeFile.txt | egrep.py "[0-9]" | line_count.py
```

Чтобы сделать то же самое в UNIX-подобной системе, потребуется немного больше работы<sup>2</sup>.

Вот аналогичный сценарий, который во входящем потоке подсчитывает слова и записывает на консоль наиболее распространенные из них:

```
# наиболее распространенные слова (most_common_words.py)
import sys
from collections import Counter

# передать число слов в качестве первого аргумента
try:
    num_words = int(sys.argv[1])
except:
    print("применение: most_common_words.py num_words")
    sys.exit(1) # ненулевой код выхода сигнализирует об ошибке

counter = Counter(word.lower() # привести слова в строчные
                  for line in sys.stdin
                  for word in line.strip().split() # разбить строку
                                                         # по пробелам
                  if word) # пропустить 'пустые' слова
```

<sup>2</sup> <http://stackoverflow.com/questions/15587877/run-a-python-script-in-terminal-without-the-python-command>.

```
for word, count in counter.most_common(num_words):
    sys.stdout.write(str(count))
    sys.stdout.write("\t")
    sys.stdout.write(word)
    sys.stdout.write("\n")
```

Затем можно сделать следующее:

```
C:\DataScience>type the_bible.txt | python most_common_words.py 10
64193 the
51380 and
34753 of
13643 to
12799 that
12560 in
10263 he
9840 shall
8987 unto
8836 for
```



Опытные UNIX-программисты, несомненно, знакомы с широким спектром встроенных инструментов командной строки, таких как `egrep`, например, которые, вероятно, будут предпочтительнее собственных, созданных с чистого листа. Тем не менее, полезно знать, что при необходимости всегда есть возможность их реализовать.

## Чтение файлов

Помимо всего прочего, можно явным образом считывать и записывать файлы непосредственно из кода. Python делает работу с файлами достаточно простой.

## Основы работы с текстовыми файлами

Первым шагом в работе с текстовым файлом является получение *файлового объекта* при помощи метода `open()`:

```
# 'r' означает только для чтения = read-only
file_for_reading = open('reading_file.txt', 'r')

# 'w' пишет в файл - сотрет файл, если он уже существует!
file_for_writing = open('writing_file.txt', 'w')

# 'a' добавляет - для добавления в конец файла
file_for_appending = open('appending_file.txt', 'a')

# не забыть закрыть файл в конце работы
file_for_writing.close()
```

Чтобы исключить случаи, когда файлы по той или иной причине не были закрыты, следует всегда пользоваться блочным оператором `with`, в конце которого они будут закрыты автоматически:

```
with open(filename, 'r') as f:
    data = function_that_gets_data_from(f)
    # в этой точке f уже был закрыт, поэтому не пытайтесь использовать его
    process(data)
```

Если требуется прочитать текстовый файл целиком, то можно выполнить перебор строк файла в цикле при помощи оператора `for`:

```
starts_with_hash = 0
```

```
with open('input.txt', 'r') as f:
    for line in file: # обратиться к каждой строке файла, используя
        if re.match("^#", line): # regex для проверки начала строки с '#'
            starts_with_hash += 1 # если да, то добавить 1 к счетчику
```

Каждая текстовая строка, полученная таким образом, заканчивается символом новой строки (или символом окончания строки), поэтому прежде чем делать с ней что-либо еще, нередко бывает необходимо удалить этот символ при помощи метода `strip()`.

Например, имеется файл с адресами электронной почты, по одному в строке, и требуется сгенерировать гистограмму доменов. Правила для корректного извлечения доменных имен несколько замысловаты (к примеру, список публичных суффиксов — <https://publicsuffix.org/>), однако в качестве первого приближения просто возьмем части адреса электронной почты, которые идут после символа `@`. (При этом такие адреса, как `joel@mail.datasciencester.com`, будут обработаны неправильно.)

```
def get_domain(email_address):
    """разбить по '@' и вернуть остаток строки"""
    return email_address.lower().split("@")[-1]
```

```
with open('email_addresses.txt', 'r') as f:
    domain_counts = Counter(get_domain(line.strip()))
    for line in f:
        if "@" in line)
```

## Файлы с разделителями

В файле с фиктивными адресами электронной почты, который был только что обработан, адреса расположены по одному в строке. Однако чаще приходится работать с файлами, где в одной строке находится большое количество данных. Значения в таких файлах часто разделены либо *запятыми*, либо *символами табуляции*. Каждая строка содержит несколько полей данных, а запятая или символ табуляции указывают, где поле заканчивается и дальше начинается новое.

Все усложняется, когда появляются поля, в которых содержатся запятые, символы табуляции и новой строки (а такие неизбежно будут). По этой причине почти

всегда будет ошибкой пытаться анализировать подобные файлы самостоятельно. Вместо этого следует использовать модуль `csv` (или библиотеку `pandas`). По техническим причинам, за которые можно свободно возложить вину на компанию Microsoft, работать с `csv`-файлами следует всегда в *двоичном* режиме, указывая флажок `b` после `r` или `w` (см. на сайте вопросов и ответов для программистов Stack Overflow<sup>3</sup>).

Если в файле нет заголовков (что, по-видимому, подразумевает, что каждая строка требуется в виде списка, и также обязывает помнить, какие данные находятся в каждом столбце), можно воспользоваться читающим объектом `csv.reader`, который выполняет навигацию по строкам, преобразуя каждую из них надлежащим образом в список.

Например, предположим, что имеется файл с курсами акций, где значения разделены символом табуляции:

6/20/2014	AAPL	90.91
6/20/2014	MSFT	41.68
6/20/2014	FB	64.5
6/19/2014	AAPL	91.86
6/19/2014	MSFT	41.51
6/19/2014	FB	64.34

Их можно обработать следующим образом:

```
import csv
```

```
with open('tab_delimited_stock_prices.txt', 'rb') as f:
    reader = csv.reader(f, delimiter='\t')
    for row in reader:
        date = row[0]
        symbol = row[1]
        closing_price = float(row[2])
        process(date, symbol, closing_price)
```

Если же файл содержит заголовки:

```
date:symbol:closing_price 6/20/2014:AAPL:90.91
6/20/2014:MSFT:41.68
6/20/2014:FB:64.5
```

то можно пропустить строку заголовков (вызвав вначале метод `reader.next()`) либо воспользоваться читающим объектом `csv.DictReader`, чтобы получить строки в виде словаря `dict` (где ключами являются заголовки):

```
with open('colon_delimited_stock_prices.txt', 'rb') as f:
    reader = csv.DictReader(f, delimiter=',')
    for row in reader:
        date = row["date"]
```

<sup>3</sup> <http://stackoverflow.com/questions/4249185/using-python-to-append-csv-files>.

```

symbol = row["symbol"]
closing_price = float(row["closing_price"])
process(date, symbol, closing_price)

```

Даже если в файле нет заголовков, можно по-прежнему использовать читающий объект DictReader, передав ему ключи в качестве аргумента fieldnames.

Запись данных с разделителями в файл выполняется аналогичным образом при помощи пишущего объекта csv.writer:

```
today_prices = { 'AAPL' : 90.91, 'MSFT' : 41.68, 'FB' : 64.5 }
```

```

with open('comma_delimited_stock_prices.txt', 'wb') as f:
    writer = csv.writer(f, delimiter=',')
    for stock, price in today_prices.items():
        writer.writerow([stock, price])

```

Пишущий объект csv.writer справится, даже если внутри поля тоже есть запятые. Собственный модуль, собранный своими силами, наверняка, на такое не будет способен. Например, если попытаться обработать самостоятельно:

```

results = [{"test1", "success", "Monday"},
           ["test2", "success, kind of", "Tuesday"],
           ["test3", "failure, kind of", "Wednesday"],
           ["test4", "failure, utter", "Thursday"]]

```

*# не делать этого!*

```

with open('bad_csv.txt', 'wb') as f:
    for row in results:
        f.write(",".join(map(str, row))) # внутри может быть много запятых!
        f.write("\n") # строка может также содержать символ новой строки!

```

то в итоге получится csv-файл, который выглядит следующим образом:

```

test1, success, Monday
test2, success, kind of, Tuesday
test3, failure, kind of, Wednesday
test4, failure, utter, Thursday

```

и который никто никогда не сможет разобрать.

## Извлечение данных из веб-ресурсов

Еще один способ получить данные — извлекать их из веб-страниц (web-scraping). Выборку самих веб-страниц, как выясняется, сделать довольно легко, а вот получить из них значимую структурированную информацию не так просто.

### Анализ кода HTML

Веб-страницы написаны на языке HTML, в котором текст (в идеале) размечен на элементы и их атрибуты:

```
<html>
  <head>
    <title>Веб-страница</title>
  </head>
  <body>
    <p id="author">Джоэл Грас</p>
    <p id="subject">Наука о данных</p>
  </body>
</html>
```

В идеальном мире, где все веб-страницы для общей пользы размечены семантически, можно было бы извлекать данные, используя правила типа "найти элемент `<p>`, чей `id` равен `subject`, и вернуть текст, который он содержит". В реальности же код на HTML обычно не только не аннотирован, он зачастую просто составлен не вполне корректно. И поэтому, чтобы в нем разобраться, придется прибегнуть к вспомогательным средствам.

Чтобы получить данные из кода HTML, применим библиотеку BeautifulSoup (<https://www.crummy.com/software/BeautifulSoup/>), которая строит дерево из разнообразных элементов, расположенных на странице, и для доступа к ним предоставляет простой интерфейс. На день написания последняя версия была BeautifulSoup 4.5.1<sup>4</sup> (`pip install beautifulsoup4`), которой как раз и воспользуемся. Кроме нее, воспользуемся библиотекой requests (<http://docs.python-requests.org/en/latest/>, `pip install requests`), которая предоставляет более удобный способ создания http-запросов, чем соответствующие средства, встроенные в Python.

Встроенный в Python синтаксический анализатор HTML слишком требователен к соблюдению правил и, как следствие, не всегда справляется с HTML, который не вполне хорошо сформирован. По этой причине будем использовать другой синтаксический анализатор, который требуется установить:

```
pip install html5lib
```

Чтобы воспользоваться библиотекой BeautifulSoup, нужно передать немного HTML-кода конструктору BeautifulSoup(). В приводимых примерах это будет результатом вызова метода requests.get():

```
from bs4 import BeautifulSoup
import requests
html = requests.get("http://www.example.com").text
soup = BeautifulSoup(html, 'html5lib')
```

И теперь можно многого добиться при помощи всего лишь нескольких простых методов.

---

<sup>4</sup> На сентябрь 2016 г. последней версией библиотеки была версия 3.5.1. Среди прилагаемых к книге материалов имеется ее установочный пакет в формате whl. Процедура установки whl-пакетов подробно описана в комментариях пользователя. — *Прим. пер.*

Как правило, мы будем работать с объектом `Tag`, который соответствует конструкциям разметки HTML или тегам, представляющим структуру HTML-страницы.

Например, следующее выражение находит первый тег `<p>` (и его содержимое):

```
first_paragraph = soup.find('p') # первый тег <p>, можно просто soup.p
```

Текстовое содержимое объекта `Tag` можно получить, воспользовавшись его свойством `text`:

```
first_paragraph_text = soup.p.text # текст первого элемента <p>
first_paragraph_words = soup.p.text.split() # слова первого элемента
```

Атрибуты тега можно извлечь, обращаясь к ним, как к словарю `dict`:

```
first_paragraph_id = soup.p['id'] # вызывает KeyError,
# если 'id' отсутствует
first_paragraph_id2 = soup.p.get('id') # возвращает None,
# если 'id' отсутствует
```

Можно получить несколько тегов сразу:

```
all_paragraphs = soup.find_all('p') # или просто soup('p')
paragraphs_with_ids = [p for p in soup('p') if p.get('id')]
```

Нередко нужно найти теги с конкретным классом `class` таблицы стилей:

```
important_paragraphs = soup('p', {'class' : 'important'})
important_paragraphs2 = soup('p', 'important')
important_paragraphs3 = [p for p in soup('p')
# if 'important' in p.get('class', [])]
```

Чтобы реализовать более детализированную логику, запросы можно комбинировать. Например, если нужно найти каждый элемент `<span>`, содержащийся внутри элемента `<div>`, то можно поступить следующим образом:

```
# элементы <span> внутри элементов <div>
# предупреждение: вернет тот же span несколько раз,
# если он находится внутри нескольких элементов div
# нужно быть смысленнее в этом случае
spans_inside_divs = [span
# для каждого <div>
# на странице
# найти каждый <span>
# внутри него
```

Всего лишь несколько функций уже позволяют достичь много. Если же все-таки требуется более сложная функциональность (или просто из любопытства), сверьтесь с документацией.

Естественно, какими бы ни были данные, как правило, их важность не обозначена в виде `class="important"`. Для обеспечения корректности данных необходимо внимательно анализировать исходный код на HTML, руководствоваться своей логикой выбора и учитывать пограничные случаи. Рассмотрим пример.



## Пример: книги об анализе данных издательства O'Reilly

Потенциальный инвестор в DataSciencester считает, что интерес к науке о данных — кратковременное явление. Чтобы доказать его неправоту, вы решаете проверить, сколько книг об анализе данных было опубликовано за длительный промежуток времени в издательстве O'Reilly. Порывшись на веб-сайте издательства, вы обнаруживаете большое число страниц из книг (и видео) о данных, до которых можно добраться через каталог, где на каждой странице выложено по 30 наименований с URL-адресами типа:

<http://shop.oreilly.com/category/browse-subjects/data.do?sortby=publicationDate&page=1>

Чтобы не выглядеть сопливом (если, конечно, вы не хотите, чтобы ваш анализатор страниц был забанен), всякий раз перед тем, как собирать данные с веб-сайта, необходимо сначала удостовериться, имеется ли там какая-нибудь политика доступа. Глядя на адрес:

<http://oreilly.com/terms/>

видно, что этому проекту вроде бы ничто не мешает. Для того чтобы оставаться порядочными людьми, следует тоже проверить файл robots.txt, который указывает поисковым роботам, как себя вести. Важными в <http://shop.oreilly.com/robots.txt> являются строки:

```
Crawl-delay: 30
Request-rate: 1/30
```

Первая сообщает, что нужно ждать 30 секунд между запросами, вторая — что следует запрашивать только одну страницу каждые 30 секунд. В сущности, строки говорят об одном и том же, только разными способами. (Есть и другие строки, которые указывают на каталоги, не подлежащие манипуляциям по сбору данных, но наш URL-адрес там не содержится, так что тут все в порядке.)



Всегда существует возможность, что издательство в какой-то момент перестроит свой сайт и сломает всю логику, показанную в этом разделе. Я, конечно, сделаю все, что смогу, чтобы помешать этому, но у меня там не мегатонны влияния. Хотя, если каждый из вас пообещает убедить всех своих знакомых купить копию этой книги...

Чтобы понять, как извлекать данные, давайте скачаем одну из этих страниц и передадим ее конструктору BeautifulSoup:

```
# не следует разбивать строку с URL-адресом таким образом,
# если не стоит задача уместить ее на странице книги
url = "http://shop.oreilly.com/category/browse-subjects/" + \
      "data.do?sortby=publicationDate&page=1"
soup = BeautifulSoup(requests.get(url).text, 'html5lib')
```

Если посмотреть на исходный код страницы (в браузере, щелкните правой кнопкой мыши и выберите команду **Просмотреть источник** или **Просмотр исходного кода**

страницы, или любой другой вариант, близкий к этому), то можно увидеть, что каждая книга (или видео) однозначно содержится в элементе ячейки таблицы `<td>` с классом `thumbtext`. Вот (сокращенный вариант) соответствующего кода на HTML для одной книги:

```
<td class="thumbtext">
  <div class="thumbcontainer">
    <div class="thumbdiv">
      <a href="/product/9781118903407.do">
        
      </a>
    </div>
  </div>
  <div class="widthchange">
    <div class="thumbheader">
      <a href="/product/9781118903407.do">Getting a Big Data Job For Dummies</a>
    </div>
    <div class="AuthorName">By Jason Williamson</div>
    <span class="directorydate">December 2014</span>
    <div style="clear:both;">
      <div id="146350">
        <span class="pricelabel">
          Ebook:

          <span class="price">&nbsp;&nbsp;&nbsp;$29.99</span>
        </span>
      </div>
    </div>
  </div>
</td>
```

Сначала хорошо бы найти все элементы тега `td` с классом `thumbtext` таблицы стилей:

```
tds = soup('td', 'thumbtext')
print(len(tds)) # 30
```

Далее необходимо отфильтровать код на HTML, удалив видео. (Будущего инвестора впечатляют только книги.) Просматривая код на HTML дальше, обнаружим, что каждая ячейка `td` содержит один или более элементов `span` с классом `pricelabel` и текстом `text`, который выглядит как `Ebook:` или `Video:`, или `Print:`. Похоже, что видеоролики содержат всего один класс `pricelabel`, чей текст `text` начинается со слова `Video` (после удаления начальных пробелов). В итоге получим следующую функцию, которая будет опознавать видео:

```
# предикативная функция - определитель видео
def is_video(td):
    """ если элемент имеет только один pricelabel,
    и текст внутри pricelabel без начальных пробелов
    начинается с 'Video', то значит это видео """
```

```

pricelabels = td('span', 'pricelabel')
return (len(pricelabels) == 1 and
        pricelabels[0].text.strip().startswith("Video"))

```

```

print(len([td for td in tds if not is_video(td)]))
# 21 (сейчас это может быть другое число)

```

Теперь можно приступить к вытягиванию данных<sup>5</sup> из элементов td. Наименование книги — это текст внутри тега <a>, который находится внутри элемента <div class="thumbheader">:

```
title = td.find("div", "thumbheader").a.text
```

Имя автора (авторов) находится в тексте элемента <div> с классом AuthorName. Имена, перед которыми стоит слово By (и которое надо удалить), перечисляются через запятую (по запятым текст надо разбить на части, после чего удалить пробелы):

```

author_name = td.find('div', 'AuthorName').text
authors = [x.strip() for x in re.sub("^By ", "", author_name).split(",")]

```

Код ISBN содержится в ссылке, которая находится в элементе <div> с классом thumbheader:

```
isbn_link = td.find("div", "thumbheader").a.get("href")
```

```

# re.match захватывает часть текста, соответствующего
# регулярному выражению в круглых скобках
isbn = re.match("/product/(.*)\.do", isbn_link).group(1)

```

Дата размещения в каталоге содержится в элементе <span class="directorydate">:

```
date = td.find("span", "directorydate").text.strip()
```

Соберем все составляющие вместе в одну функцию:

```

# информация о книге
def book_info(td):
    """Имя на входе объект Tag для <td> библиотеки BeautifulSoup, который
    обозначает книгу, извлечь описание книги и вернуть словарь dict"""

    title = td.find("div", "thumbheader").a.text
    by_author = td.find('div', 'AuthorName').text
    authors = [x.strip() for x in re.sub("^By ", "", by_author).split(",")]
    isbn_link = td.find("div", "thumbheader").a.get("href")
    isbn = re.match("/product/(.*)\.do", isbn_link).groups()[0]
    date = td.find("span", "directorydate").text.strip()

    return {
        "title" : title,

```

<sup>5</sup> Вытягивание данных (data pull) — технология поиска пользователем информации в базах данных и в веб-ресурсах. — Прим. пер.

```

"authors" : authors,
"isbn" : isbn,
"date" : date
}

```

Теперь все готово для извлечения данных со страниц:

```

from bs4 import BeautifulSoup
import requests
from time import sleep

base_url = "http://shop.oreilly.com/category/browse-subjects/" + \
           "data.do?sortBy=publicationDate&page="

books = []

NUM_PAGES = 40 # данные на конец 2015 (на момент написания было 31)

for page_num in range(1, NUM_PAGES + 1):
    print("сбор данных на странице", page_num, ", всего найдено", len(books))
    url = base_url + str(page_num)
    soup = BeautifulSoup(requests.get(url).text, 'html5lib')

    for td in soup('td', 'thumbtext'):
        if not is_video(td): books.append(book_info(td))

# порядочный человек должен уважать robots.txt!
sleep(30)

```



Такого рода извлечение данных из HTML больше походит на творчество, чем на науку о данных. Существует бесчисленное множество других схем поиска наименований книг, фильмов и других источников, которые тоже будут работать хорошо.

Теперь, когда данные извлечены, можно построить график количества книг, издававшихся в каждый год (рис. 9.1):

```

# получить год
def get_year(book):
    """book["date"] возвращает 'December 2015', поэтому надо
    разбить строку по пробелу и затем взять вторую часть"""
    return int(book["date"].split()[1])

# число лет
# 2015 год - последний полный год данных (на момент тестирования сценария)
year_counts = Counter(get_year(book) for book in books
                       if get_year(book) <= 2015)

import matplotlib.pyplot as plt
years = sorted(year_counts)
book_counts = [year_counts[year] for year in years]

```

```
plt.plot(years, book_counts)
plt.ylabel("Число книг по анализу данных")
plt.axis([2000,2016,0,260]) # установить оси
plt.title("Данные большие!")
plt.show()
```

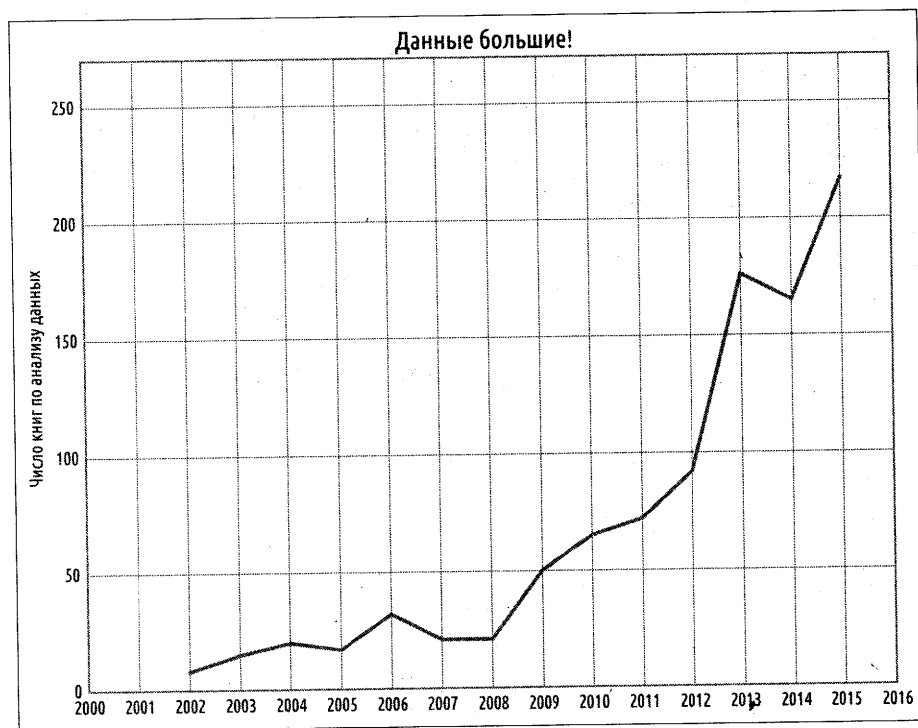


Рис. 9.1. Число книг по анализу данных, издаваемых ежегодно

Потенциальный инвестор смотрит на диаграмму и решает, что 2015 год был пиковым для изданий этой направленности.

## Использование программных интерфейсов

Многие веб-сайты и веб-сервисы предоставляют интерфейсы программирования приложений (Applications Programming Interface, API) или просто программных интерфейсов, которые позволяют явным образом запрашивать данные в структурированном формате. И это в значительной мере избавляет от необходимости извлекать данные самому!

### Формат JSON (и XML)

Поскольку протокол HTTP служит для передачи текста, то данные, которые запрашиваются через веб-API, должны быть преобразованы в строковый формат или

сериализованы. Нередко при сериализации используется объектная нотация языка JavaScript (JavaScript Object Notation, JSON). Объекты в JavaScript очень похожи на словари dict в языке Python, что облегчает интерпретацию их строковых представлений:

```
{ "title" : "Data Science Book",
  "author" : "Joel Grus",
  "publicationYear" : 2014,
  "topics" : [ "data", "science", "data science" ] }
```

Текст в формате JSON анализируется при помощи модуля json, в частности, посредством его функции loads, которая преобразует (десериализует) строку, представляющую объект JSON, в объект на Python:

```
import json
serialized = """{ "title" : "Data Science Book",
                  "author" : "Joel Grus",
                  "publicationYear" : 2014,
                  "topics" : [ "data", "science", "data science" ] }"""
```

```
# преобразовать строку в формате JSON в питоновский словарь dict
deserialized = json.loads(serialized)
if "data science" in deserialized["topics"]:
    print(deserialized)
```

Иногда поставщик API вас ненавидит и предоставляет ответы только в формате XML:

```
<Book>
  <Title>Data Science Book</Title>
  <Author>Joel Grus</Author>
  <PublicationYear>2014</PublicationYear>
  <Topics>
    <Topic>data</Topic>
    <Topic>science</Topic>
    <Topic>data science</Topic>
  </Topics>
</Book>
```

Чтобы получить данные из XML, можно воспользоваться библиотекой BeautifulSoup аналогично тому, как она использовалась для получения данных из HTML. Сверьтесь с документацией касательно подробностей.

## Использование непроверенного API

Большинство интерфейсов программирования приложений (API) в наши дни требуют от пользователя сначала аутентифицировать себя прежде, чем начать ими пользоваться. Хотя никто не порицает их за такую политику, но так или иначе она создает много излишних шаблонных процедур, которые затуманивают непосред-

венный контакт. Поэтому сначала обратимся к программному интерфейсу GitHub<sup>6</sup>, с помощью которого можно делать некоторые простые вещи без проверки подлинности:

```
import requests, json
endpoint = "https://api.github.com/users/joelgrus/repos" # конечная точка

repos = json.loads(requests.get(endpoint).text)
```

В этом месте переменная `repos` представляет собой список из словарей `dict`, каждый из которых обозначает публичное хранилище в моем аккаунте на GitHub. (Можете подставить свое имя пользователя и получить собственные данные из хранилища. Ведь у вас уже есть аккаунт на GitHub?)

Эти данные можно использовать, например, чтобы выяснить, в какие месяцы и по каким дням недели я чаще всего создаю хранилища. Единственная проблема заключается в том, что даты в ответе представлены строкой символов (Юникода):

```
u'created_at': u'2013-07-05T02:02:28Z'
```

Python поставляется без сносного анализатора дат, поэтому придется установить свой:

```
pip install python-dateutil
```

в котором, скорее всего, понадобится только функция `dateutil.parser.parse`:

```
from dateutil.parser import parse

dates = [parse(repo["created_at"]) for repo in repos] # список дат
month_counts = Counter(date.month for date in dates) # число месяцев
weekday_counts = Counter(date.weekday() for date in dates) # число
# будних дней
```

Таким же образом можно получить языки программирования моих последних пяти хранилищ:

```
last_5_repositories = sorted(repos, # последние 5 хранилищ
                             key=lambda r: r["created_at"],
                             reverse=True)[:5]

last_5_languages = [repo["language"] # последние 5 языков
                    for repo in last_5_repositories]
```

Как правило, работа с API будет выполняться не на низком уровне в виде выполнения запросов и самостоятельного разбора ответов. Одно из преимуществ использования языка Python заключается в том, что кто-то уже разработал библиотеку для практически любого API, доступ к которому требуется получить. Когда библиотеки

<sup>6</sup> Github — один из самых популярных сервисов для создания, распространения и совместной работы с программным обеспечением (<https://developer.github.com/v3/>). — Прим. пер.

выполнены сносно, они могут сэкономить уйму неприятностей при выяснении проблемных подробностей доступа к API. (Когда же они выполнены кое-как или выясняется, что они основаны на несуществующей версии соответствующих API, они на самом деле становятся причиной головной боли.)

Так или иначе, периодически возникает необходимость разворачивать собственную библиотеку доступа к API (или, что более вероятно, копаться в причинах, почему чужая библиотека не работает, как надо), поэтому хорошо бы познакомиться с некоторыми подробностями.

## Поиск API

Если необходимо получить данные с определенного сайта, то все подробности следует искать в разделе для разработчиков или разделе API сайта. Помимо этого, можно попытаться отыскать библиотеку в Интернете по запросу "python \_\_ api". Вот лишь некоторые из них: API на Python для веб-сервиса Rotten Tomatoes, предоставляющего статистику и обзоры кинофильмов, ряд оберток (или надстроек)<sup>7</sup>, таких как обертка API на Python для веб-сайта Klout, предоставляющего аналитику социальных связей на основе ведущих мировых соцсетей, обертка API для платформы Yelp для интегрирования транзакционных услуг, предлагаемых сторонними разработчиками местным предприятиям, обертка API для веб-сервиса IMDb для выполнения поиска по кинематографической базе данных Imdb и т. д.

Если нужны списки API с обертками на Python, то вот три каталога: API-надстройки на Python (<https://github.com/realpython/list-of-python-api-wrappers>), Python API и Python для начинающих (<http://www.pythonforbeginners.com/development/list-of-python-apis/>).

Если нужен более широкий каталог веб-API для создания служб HTTP (не обязательно с обертками на Python), то хорошим ресурсом будет Programmable Web (<http://www.programmableweb.com/>), где есть огромный каталог, распределенный по категориям API.

И если после всего этого вы так ничего и не нашли из того, что вам нужно, то всегда можно обратиться к извлечению данных с веб-страниц как к последнему прибежищу аналитика данных.

## Пример: использование интерфейсов Twitter API

Twitter — это фантастический источник данных для работы. Его можно использовать для получения новостей в реальном времени, для измерения реакции на текущие события, для поиска ссылок, связанных с конкретными темами, практически

---

<sup>7</sup> Обертка API — это промежуточный слой между средой разработки на Python и интерфейсом программирования приложений (API), предоставляемым веб-сервисом, распределенной платформой или библиотекой общего пользования. — *Прим. пер.*



для всего, что можно представить себе, пока имеется доступ к его данным, который можно получить через его API.

Для взаимодействия с интерфейсами программирования приложений, предоставленными разработчикам социальной сетью Twitter, будем использовать библиотеку Twython (<https://github.com/ryanmcgrath/twython>, `pip install twython`). Для работы с соцсетью Twitter на языке Python имеется целый ряд других библиотек, однако работа именно с этой библиотекой у меня удалась лучше всего. Вы можете исследовать и другие!

## Получение учетных данных

Для того чтобы воспользоваться интерфейсами программирования приложений (API) соцсети Twitter, нужно получить некоторые полномочия (для которых необходимо иметь учетную запись Twitter, и которую так или иначе следует иметь, благодаря чему можно стать частью живого и дружелюбного сообщества Twitter #datascience). Как и все остальные инструкции, касающиеся веб-сайтов, которые я не контролирую, в определенный момент они могут устареть, но все же есть надежда, что они будут актуальными в течение некоторого времени. (Хотя они уже менялись, по крайней мере один раз, пока я писал эту книгу, так что удачи!)

Итак:

1. Перейдите на <https://apps.twitter.com/>.
2. Если вы не вошли на сайт, щелкните на ссылке **Sign in** и введите имя пользователя Twitter и пароль.
3. Нажмите кнопку **Create New App** (Создать новое приложение).
4. Присвойте ему имя (например, "наука о данных"), дайте описание и укажите любой URL-адрес в качестве веб-сайта (не важно, какой).
5. Согласитесь с условиями использования и нажмите кнопку **Create Your Twitter Application** (Создайте ваше приложение Twitter).
6. Запишите ключ потребителя (**Consumer Key**) и секрет потребителя (**Consumer Secret**).
7. Нажмите кнопку **Create my access token** (Создать мой маркер доступа).
8. Запишите маркер доступа и секретную часть маркера доступа (может понадобиться обновить страницу).

Ключ потребителя и секрет потребителя сообщают в Twitter, какое приложение обращается к его интерфейсам API, в то время как маркер доступа и секретная часть маркера доступа говорят о том, кто это делает. Если вы когда-либо использовали собственный аккаунт Twitter для входа на какой-нибудь другой сайт, то страница со ссылкой **Нажмите, чтобы разрешить** как раз генерирует маркер доступа для этого сайта, который используется для того, чтобы убедить Twitter, что это были вы (или, по крайней мере, кто-то, действующий от вашего имени). Поскольку нам не нужна такая открытая функциональность ("любой может войти"), обойдемся статически созданными маркером доступа и секретной частью маркера доступа.



Ключ/секрет потребителя и ключ/секрет маркера доступа должны рассматриваться как *пароли*. Вы не должны делиться ими, публиковать их в своем блоге и передавать их в ваше публичное хранилище на GitHub. Одно из простых решений — хранить их в файле `credentials.json`, который не передается, а в своем коде использовать метод `json.loads` для их извлечения.

## Использование Twython

Сначала обратимся к API поиска в соцсети Twitter Search API<sup>8</sup> (<https://dev.twitter.com/rest/reference/get/search/tweets>), который требует только ключ и секрет потребителя (ключ/секрет маркера доступа не требуются):

```
from twython import Twython
```

```
twitter = Twython(CONSUMER_KEY, CONSUMER_SECRET)
```

```
# искать твиты, содержащие фразу "data science"
for status in twitter.search(q="data science")["statuses"]:
    user = status["user"]["screen_name"].encode('utf-8')
    text = status["text"].encode('utf-8')
    print(user, ":", text)
print()
```



Метод `.encode("utf-8")` нужен, чтобы учитывать случаи, когда твиты содержат символы Юникода, с которыми оператор `print` не справляется. (Если его убрать, то, скорее всего, будет получена ошибка `UnicodeEncodeError`.)

В определенный момент вашей деятельности в области аналитики данных вы почти наверняка столкнетесь с серьезными трудностями, связанными с Юникодом, что заставит вас обратиться к документации по Python<sup>9</sup> или же нехотя начать использовать Python 3, который обходится с текстом в Юникоде гораздо дружелюбнее.

Если выполнить этот сценарий, то можно получить несколько твитов, наподобие следующих:

```
haithemnyс: Data scientists with the technical savvy & analytical chops to
derive meaning from big data are in demand. http://t.co/HsF9Q0dShP
```

```
RPubсRecent: Data Science http://t.co/6hсHuz2PHM
```

```
spleonard1: Using #dplyr in #R to work through a procrastinated assignment
for @rdpeng in @coursera data science specialization. So easy and Awesome.
```

<sup>8</sup> API поиска в соцсети Twitter (Twitter Search API) — это компонент API социальной сети Twitter, построенный на основе архитектуры взаимодействия компонентов распределенного приложения в сети (REST), который позволяет делать запросы относительно показателей последних или популярных твитов и ведет себя аналогично функциям поиска доступных для мобильных или веб-клиентов, таких как поиск в Twitter.com. API поиска в соцсети Twitter выполняет поиск по выборке недавних твитов, опубликованных за последние 7 дней. — *Прим. пер.*

<sup>9</sup> <https://docs.python.org/2/howto/unicode.html>.

Эта информация не представляет особого интереса, главным образом потому, что API поиска в соцсети Twitter Search API всего лишь показывает несколько последних результатов, которые он считает нужным показать. Когда вы решаете задачи, связанные с наукой о данных, чаще требуется гораздо больше твитов. Для таких целей служит Streaming API<sup>10</sup>. Этот интерфейс позволяет подключаться к огромному потоку сообщений Twitter. Чтобы им воспользоваться, необходимо аутентифицироваться с помощью личного маркера доступа.

Чтобы получить доступ к Streaming API при помощи Twython, необходимо определить класс, который наследует у класса TwythonStreamer и переопределяет его метод `on_success()` (и, возможно, его метод `on_error()`):

```
from twython import TwythonStreamer

# добавлять данные в глобальную переменную - это пример плохого стиля,
# но он намного упрощает пример
tweets = []

class MyStreamer(TwythonStreamer):
    """ наш собственный подкласс класса TwythonStreamer, который
    определяет, как взаимодействовать с потоком """

    def on_success(self, data):
        """ что делать, когда twitter присылает данные?
        здесь данные будут в виде словаря dict, представляющего твит """

        # нужно собирать твиты только на английском
        if data['lang'] == 'en':
            tweets.append(data)
            print("received tweet #", len(tweets))

        # остановиться, когда собрано достаточно
        if len(tweets) >= 1000:
            self.disconnect()

    def on_error(self, status_code, data):
        print(status_code, data)
        self.disconnect()
```

Подкласс `MyStreamer` подключается к потоку Twitter и ждет, когда Twitter подаст данные. Каждый раз, когда он получает некоторые данные (здесь, твит представлен как объект языка Python), он передает этот твит в метод `on_success()`, который добавляет его к списку полученных твитов `tweets` при условии, что они на английском языке, и в конце, после сбора 1000 твитов, он отсоединяет стример.

<sup>10</sup> Поточковые интерфейсы API соцсети Twitter Streaming API — это ряд компонентов API социальной сети Twitter, построенных на основе архитектуры REST, которые позволяют разработчикам получать доступ к глобальному потоку твит-данных с низкой задержкой доступа. — Прим. пер.

Осталось только инициализировать его и запустить:

```
# запуск потока
stream = MyStreamer(CONSUMER_KEY, CONSUMER_SECRET,
                    ACCESS_TOKEN, ACCESS_TOKEN_SECRET)

# начинает потреблять публичные статусы,
# которые содержат ключевое слово 'data'
stream.statuses.filter(track='data')

# в случае, если напротив нужно начать потреблять ВСЕ публичные статусы
# stream.statuses.sample()
```

Этот сценарий будет выполняться до тех пор, пока он не соберет 1000 твитов (или пока не встретит ошибку), и остановится, после чего можно приступить к анализу твитов. Например, наиболее распространенные хештеги можно найти следующим образом:

```
# ведущие хештеги
top_hashtags = Counter(hashtag['text'].lower()
                       for tweet in tweets
                       for hashtag in tweet["entities"]["hashtags"])

print(top_hashtags.most_common(5))
```

Каждый твит содержит большое количество данных. За подробностями можно обратиться к окружающим либо покопаться в документации Twitter API<sup>11</sup>.



В реальном проекте, вероятно, вместо того чтобы полагаться на список, хранящийся в оперативной памяти, как на структуру данных для твитов, следует сохранить их в файл или базу данных, чтобы иметь их всегда под рукой.

## Для дальнейшего изучения

- ◆ **pandas** (<http://pandas.pydata.org/>) — основная библиотека, используемая в науке о данных для работы с данными (в частности, для импорта данных).
- ◆ **Scrapy** (<http://scrapy.org/>) — более полнофункциональная библиотека для построения сложных веб-анализаторов, которые обладают такой функциональностью, как переход по неизвестным ссылкам.

<sup>11</sup> <https://dev.twitter.com/overview/api/tweets><http://pandas.pydata.org/>.

# Обработка данных

Эксперты часто обладают, скорее, данными, чем здравым смыслом.

Колин Пауэлл<sup>1</sup>

Обработка данных имеет как научно-теоретический, так и практический аспекты. До сих пор главным образом обсуждалась научно-теоретическая сторона работы с данными. В этой главе будут рассмотрены некоторые профессиональные приемы, применяемые на практике.

## Исследование данных

После того как определены вопросы, на которые необходимо получить ответ, и на руках имеются некоторые данные, может возникнуть соблазн сразу же погрузиться в работу, начав строить модели и получать ответы. Не следует поддаваться этому искушению. Первым шагом должно стать, прежде всего, *исследование* данных.

## Исследование одномерных данных

В простейшем случае может иметься одномерный набор данных, который представляет собой всего лишь ряд чисел. Например, это может быть среднее число минут, которые каждый пользователь проводит на сайте каждый день, число просмотров каждого видеоролика из коллекции, посвященной анализу данных, или число страниц в каждой книге, посвященной науке о данных, из библиотеки.

Очевидно, первым делом необходимо вычислить несколько сводных статистик. К примеру, можно выяснить количество точек данных, их минимальное, максимальное значения, среднее значение и стандартное отклонение.

Впрочем, эти показатели не всегда способны помочь разобраться в данных. На следующем этапе надо построить гистограмму, в которой данные разбиты на дискретные *интервалы* (бакеты), и подсчитать число точек данных, попадающих в каждый из них:

```
# привести точку данных к номеру интервала
def bucketize(point, bucket_size):
```

---

<sup>1</sup> Колин Лютер Пауэлл (1937) — генерал Вооруженных сил США, госсекретарь в период первого срока президентства Дж. Буша-младшего. Известен в связи с фальсификацией в Совете Безопасности ООН 5 февраля 2003 г., когда он демонстрировал ампулу, якобы, с бактериями сибирской язвы, чтобы подтвердить факт наличия биологического оружия в Ираке. — *Прим. пер.*

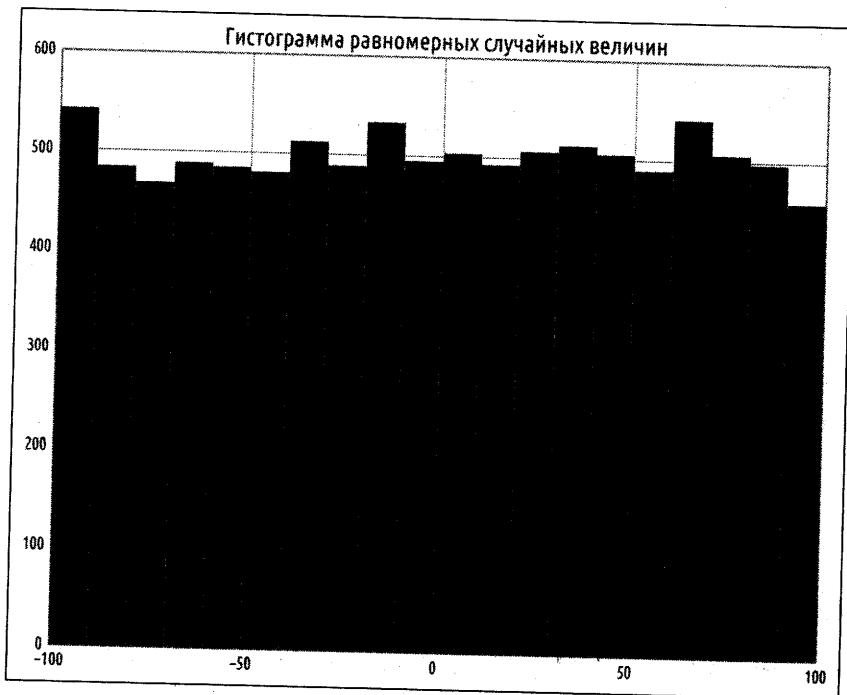


Рис. 10.1. Гистограмма равномерно распределенных случайных величин *uniform*



Рис. 10.2. Гистограмма нормально распределенных случайных величин *normal*

```

"""округлить точку до следующего наименьшего кратного
размера интервала bucket_size"""
    return bucket_size * math.floor(point / bucket_size)

# подготовить гистограмму
def make_histogram(points, bucket_size):
    """сгруппировать точки и подсчитать количество в интервале"""
    return Counter(bucketize(point, bucket_size) for point in points)

# изобразить гистограмму
def plot_histogram(points, bucket_size, title=""):
    histogram = make_histogram(points, bucket_size)
    plt.bar(histogram.keys(), histogram.values(), width=bucket_size)
    plt.title(title)
    plt.show()

```

Например, рассмотрим следующие два набора данных:

```

random.seed(0)

# равномерно распределенные между -100 и 100
uniform = [200 * random.random() - 100 for _ in range(10000)]

# нормально распределенные с нулевым средним, стандартным отклонением 57
normal = [57 * inverse_normal_cdf(random.random())
          for _ in range(10000)]

```

У обоих средние значения близки к 0, а стандартные отклонения — к 58. Тем не менее, они имеют совершенно разные распределения. На рис. 10.1 показано равномерно распределенный набор данных `uniform`:

```
plot_histogram(uniform, 10, "Гистограмма равномерных величин")
```

а на рис. 10.2 — нормально распределенный набор данных `normal`:

```
plot_histogram(normal, 10, "Гистограмма нормальных величин")
```

В данном случае у обоих распределений довольно разные показатели `max` и `min`, но даже этих сведений недостаточно, чтобы понять *степень* их различия.

## Двумерные данные

Теперь в нашем распоряжении двумерный набор данных. Например, кроме ежедневного числа минут может еще учитываться годовой стаж работы в области анализа данных. Разумеется, можно разобраться в каждой размерности по отдельности. А можно, скажем, раскидать данные на точечной диаграмме.

Для примера рассмотрим еще один фиктивный набор данных:

```

# случайная выборка из нормального распределения
def random_normal():

```

```
"""возвращает случайную выборку из стандартного
нормального распределения"""
```

```
return inverse_normal_cdf(random.random())
```

```
xs = [random_normal() for _ in range(1000)]
ys1 = [x + random_normal() / 2 for x in xs]
ys2 = [-x + random_normal() / 2 for x in xs]
```

Если бы нужно было выполнить функцию `plot_histogram` со списками `ys1` и `ys2` в качестве аргументов, то получились бы очень похожие диаграммы (действительно, обе нормально распределены с одинаковым средним значением и стандартным отклонением).

Но у обоих, как видно из рис. 10.3, имеются большие различия в совместном распределении по `xs`:

```
plt.scatter(xs, ys1, marker='.', color='black', label='ys1')
plt.scatter(xs, ys2, marker='.', color='gray', label='ys2')
plt.xlabel('xs')
plt.ylabel('ys')
plt.legend(loc=9)
plt.title("Очень разное совместное распределение")
plt.show()
```

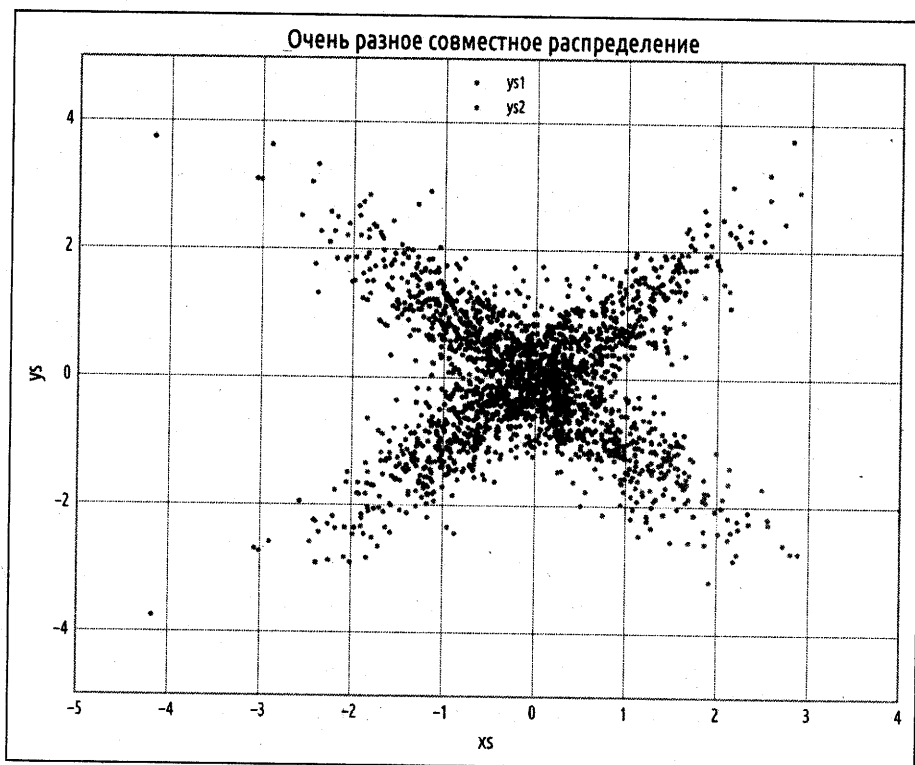


Рис. 10.3. Точечная диаграмма двух разных наборов данных `ys`





```

# затем спрятать осевые метки за исключением левой и нижней диаграмм
if i < num_columns - 1: ax[i][j].xaxis.set_visible(False)
if j > 0: ax[i][j].yaxis.set_visible(False)

# настроить правую нижнюю и левую верхнюю осевые метки,
# которые некорректны, потому что на их диаграммах выводится только текст
ax[-1][-1].set_xlim(ax[0][-1].get_xlim())
ax[0][0].set_ylim(ax[0][1].get_ylim())

plt.show()

```

Рассматривая точечные диаграммы, можно увидеть, что серия 1 очень отрицательно коррелирует с серией 0, серия 2 положительно коррелирует с серией 1, а серия 3 принимает только значения 0 и 6, причем 0 соответствует малым значениям серии 2, а 6 соответствует большим значениям.

Точечная матрица позволяет быстро получить общее представление о том, какие переменные коррелируют (если, конечно, не сидеть часами, пытаясь точно настроить matplotlib, чтобы изобразить диаграмму именно так, как хочется, и тогда это не быстрый способ).

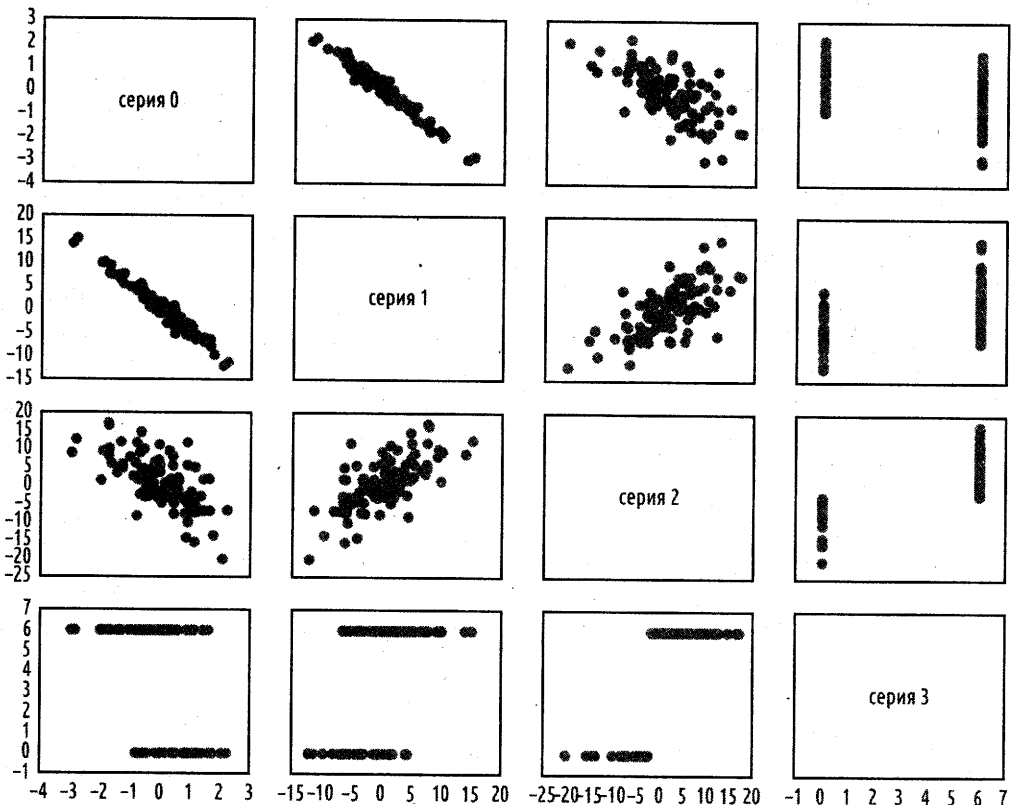


Рис. 10.4. Точечная матрица

## Очистка и форматирование

В реальном мире приходится иметь дело с *грязными* данными. Прежде чем данные можно будет использовать, нередко их необходимо подвергнуть предварительной обработке. Подобные примеры были показаны в *главе 9*. Прежде чем начать пользоваться числами, сначала надо преобразовать строки в вещественные числа *float* или целые числа *int*. Ранее это делалось непосредственно перед использованием данных:

```
# цена на момент закрытия биржи
closing_price = float(row[2])
```

Однако, вероятно, менее подверженный ошибкам способ — преобразовывать данные на входе. Это можно реализовать, создав функцию-обертку для читающего объекта *csv.reader*. В нее будет передаваться список анализаторов, каждый из которых задает способ преобразования одного из столбцов, где *None* означает, что с этим столбцом ничего делать не надо:

```
# разобрать строку табличных данных при помощи анализатора parsers(value)
def parse_row(input_row, parsers):
    """применить соответствующий анализатор из заданного списка (некоторые
    могут быть пустыми (None)) к каждому элементу input_row"""

    return [parser(value) if parser is not None else value
            for value, parser in zip(input_row, parsers)]

# разобрать строки табличных данных
def parse_rows_with(reader, parsers):
    """обертывает объект reader, применяя анализаторы к каждой строке"""
    for row in reader:
        yield parse_row(row, parsers)
```

Как быть в тех случаях, когда данные плохие? Или когда встретилось "вещественное" значение, которое на самом деле не является числом? В таких случаях обычно лучше получить специальное значение *None*, чем столкнуться со сбоем в работе программы. С этим может справиться вспомогательная функция:

```
# попытаться разобрать или вернуть None
def try_or_none(f):
    """обертывает функцию f, возвращая None, если f вызывает исключение;
    подразумевается, что f - функция одного входящего аргумента"""
    def f_or_none(x):
        try: return f(x)
        except: return None
    return f_or_none
```

Затем можно переписать функцию анализа строки данных *parse\_row*, чтобы использовать ее в дальнейшем:

*# разобрать строку табличных данных (с функцией-оберткой)*

```
def parse_row(input_row, parsers):
    return [try_or_none(parser)(value) if parser is not None else value
            for value, parser in zip(input_row, parsers)]
```

Например, если список курсов акций с разделителями-запятыми содержит плохие данные:

```
6/20/2014,AAPL,90.91
6/20/2014,MSFT,41.68
6/20/2014,FB,64.5
6/19/2014,AAPL,91.86
6/19/2014,MSFT,n/a
6/19/2014,FB,64.34
```

то теперь этот список можно прочитать и преобразовать всего за один шаг:

```
import dateutil.parser
data = []

with open("comma_delimited_stock_prices.csv", "rb") as f:
    reader = csv.reader(f)
    for line in parse_rows_with(reader, [dateutil.parser.parse, None, float]):
        data.append(line)
```

Далее остается лишь проверить, имеются ли строки с None:

```
for row in data:
    if any(x is None for x in row):
        print(row)
```

и решить, что с ними делать. (Вообще говоря, есть три варианта: удалить, попытаться исправить плохие/отсутствующие данные, обратившись к источнику, либо не делать ничего, скрестив пальцы.)

Аналогичные функции-помощники можно создать для читающего объекта `csv.DictReader`, предназначенного для работы со словарями. В этом случае, вероятно, нужно передать словарь `dict` с анализаторами, именованными по названию поля. Например:

```
# попытаться преобразовать поле в число
def try_parse_field(field_name, value, parser_dict):
    """ попытаться преобразовать значение при помощи соответствующей
    функции из словаря parser_dict с анализаторами """
    parser = parser_dict.get(field_name) # None, если нет такой записи
    if parser is not None:
        return try_or_none(parser)(value)
    else:
        return value
```

*# проанализировать (разобрать) словарь с данными (с функцией-оберткой)*

```
def parse_dict(input_dict, parser_dict):
```

```
return { field_name : try_parse_field(field_name, value, parser_dict)
        for field_name, value in input_dict.iteritems() }
```

Следующим шагом будет проверка данных на наличие выбросов, используя приемы из *разд. "Исследование данных"* текущей главы, или решений, исходя из конкретной ситуации. Например, как уже можно было заметить, одна из дат в файле с ценами на акции была 3014 год. Это значение не вызовет (с неизбежностью) ошибку, но оно в корне неверно, и если его не идентифицировать, то в итоге можно получить нелепый результат. В реальных массивах данных могут отсутствовать десятичные разделители, иметься лишние нули, опечатки и другие бесчисленные проблемы, которые придется отлавливать. (Эта работа может и не быть официально прописана в должностной инструкции, но кто же еще этим займется?)

## Управление данными

Одним из самых важных профессиональных навыков аналитика данных является умение *управлять данными*. Это, скорее, общий подход, чем конкретный набор приемов, поэтому пройдемся по некоторым примерам, чтобы почувствовать его особенности.

Представим, что мы работаем со словарями `dict`, содержащими курсы акций в следующем формате:

```
data = [
    {'closing_price': 102.06,
     'date': datetime.datetime(2014, 8, 29, 0, 0),
     'symbol': 'AAPL'},
    # ...
]
```

Концептуально представим их в виде строк (как в электронной таблице).

Начнем задавать вопросы об этих данных, попутно пытаюсь обнаружить закономерности или шаблоны в том, что делается, и помечая для себя некоторые инструменты, которые облегчают управление данными.

Например, предположим, что нужно узнать самую высокую цену на момент закрытия торгов для акций AAPL<sup>2</sup>. Разобьем задачу на конкретные шаги:

1. Ограничиться строками с символом AAPL.
2. Извлечь в каждой строке цену закрытия `closing_price`.
3. Взять из этих цен максимум `max`.

Все три шага можно выполнить одновременно при помощи генератора последовательности:

```
# максимальная цена на акции AAPL
max_aapl_price = max(row["closing_price"]
                    for row in data
                    if row["symbol"] == "AAPL")
```

---

<sup>2</sup> Акции Apple Inc. — *Ред.*

В более общем плане, хотелось бы получать сведения о самых высоких показателях цены на момент закрытия для каждой ценной бумаги из набора данных. Вот один из способов, как это сделать:

1. Сгруппировать все строки с одинаковым символом `symbol`.
2. В пределах каждой группы, проделать то же, что было описано раньше:

```
# сгруппировать все строки на основе символа
by_symbol = defaultdict(list)
for row in data:
    by_symbol[row["symbol"]].append(row)

# использовать генератор последовательности в виде словаря,
# чтобы найти максимум для каждого символа
max_price_by_symbol = { symbol : max(row["closing_price"]
                                     for row in grouped_rows)
                       for symbol, grouped_rows in by_symbol.iteritems() }
```

Уже здесь имеются некие шаблоны. В обоих случаях нужно вытянуть значение цены на момент закрытия `closing_price` из каждого словаря. Поэтому создадим функцию для получения значения поля из словаря `dict`, а другую — для извлечения значения того же поля из набора словарей:

```
# вернуть безымянную функцию, которая по полю field_name
# возвращает значение из словаря row
def picker(field_name):
    """возвращает функцию, которая извлекает поле из dict"""
    return lambda row: row[field_name]

# получить список значений из словаря по полю field_name
def pluck(field_name, rows):
    """преобразует список словарей dict в список значений
    согласно field_name"""
    return map(picker(field_name), rows)
```

Кроме того, можно создать функцию, которая группирует строки по результату группирующей функции `grouper`, применяя для каждой группы в случае необходимости один из преобразователей значений `value_transform`:

```
# группировать по
def group_by(grouper, rows, value_transform=None):
    # ключ - это результат вычисления grouper, значение - это список строк
    grouped = defaultdict(list)
    for row in rows:
        grouped[grouper(row)].append(row)

    if value_transform is None:
        return grouped
```

```

else:
    return { key : value_transform(rows)
            for key, rows in grouped.iteritems() }

```

Это позволит довольно легко переписать предыдущие примеры. Например:

```

max_price_by_symbol = group_by(picker("symbol"),
                               data,
                               lambda rows: max(pluck("closing_price", rows)))

```

Теперь можно задавать более сложные вопросы, а именно, каковы максимальные и минимальные однодневные процентные изменения в наборе данных. Процентное изменение  $\text{price\_today} / \text{price\_yesterday} - 1$  подразумевает, что необходимо каким-то образом увязать сегодняшнюю цену со вчерашней. Один из способов состоит в том, чтобы сгруппировать цены по символу и затем в пределах каждой группы:

1. Упорядочить цены по дате.
2. Получить пары цен в формате (предыдущие, текущие) при помощи встроенной функции `zip`.
3. Преобразовать пары в новые строки с процентными изменениями.

Начнем с написания функции, которая выполняет всю работу в пределах одной группы:

```

# процентное изменение цены
def percent_price_change(yesterday, today):
    return today["closing_price"] / yesterday["closing_price"] - 1

# изменения день ото дня
def day_over_day_changes(grouped_rows):
    # сортировать строки по дате
    ordered = sorted(grouped_rows, key=picker("date"))

    # объединить со смещением, чтобы получить последовательные пары
    return [{ "symbol" : today["symbol"],
              "date" : today["date"],
              "change" : percent_price_change(yesterday, today) }
            for yesterday, today in zip(ordered, ordered[1:])]

```

Затем последнюю функцию можно использовать непосредственно в качестве преобразователя `value_transform` в функции `group_by` группировки по полю:

```

# ключ - это символ,
# значение - это список из словарей dict с изменениями "change"
changes_by_symbol = group_by(picker("symbol"), data, day_over_day_changes)

# собрать все словари с изменениями "change" в один большой список
all_changes = [change
                for changes in changes_by_symbol.values()
                for change in changes]

```

На данном этапе можно легко найти максимум и минимум:

```
max(all_changes, key=picker("change"))
# {'change': 0.3283582089552237,
#  'date': datetime.datetime(1997, 8, 6, 0, 0),
#  'symbol': 'AAPL'}
# например, см. http://news.cnet.com/2100-1001-202143.html
```

```
min(all_changes, key=picker("change"))
# {'change': -0.5193370165745856,
#  'date': datetime.datetime(2000, 9, 29, 0, 0),
#  'symbol': 'AAPL'}
# например, см. http://money.cnn.com/2000/09/29/markets/techwrap/
```

Теперь можно воспользоваться новым набором данных обо всех изменениях `all_changes`, чтобы найти месяц, в котором лучше всего инвестировать в технологические акции. Сначала сгруппируем изменения по месяцам, затем вычислим суммарное изменение в пределах каждой группы.

Еще раз напишем соответствующий преобразователь значений `value_transform` и затем применим функцию `group_by` группировки по полю:

```
# объединить процентные изменения
# для этого добавим 1 к каждому, умножим и отнимем 1
# например, при объединении +10% и -20% суммарное изменение будет
# (1 + 10%) * (1 - 20%) - 1 = 1.1 * .8 - 1 = -12%
def combine_pct_changes(pct_change1, pct_change2):
    return (1 + pct_change1) * (1 + pct_change2) - 1

# суммарное изменение
def overall_change(changes):
    return reduce(combine_pct_changes, pluck("change", changes))

# суммарное ежемесячное изменение
overall_change_by_month = group_by(lambda row: row['date'].month,
                                   all_changes,
                                   overall_change)
```

Аналогичные приемы управления данными будут применяться на протяжении всей книги, обычно, не привлекая к ним особого внимания.

## Шкалирование

Многие методы чувствительны к *шкале* данных. Например, пусть дан массив сведений о росте и весе сотен аналитиков данных, и нужно определить *кластеры* размеров тела.

Логично представить кластеры в виде точек, расположенных рядом друг с другом, а значит, потребуется некоторое представление о расстоянии между точками. У нас



уже есть функция `distance`, вычисляющая евклидово расстояние, поэтому совершенно естественно рассматривать пары в формате (рост, вес) как точки в двумерном пространстве. Посмотрим на данные нескольких человек, представленные в табл. 10.1.

Таблица 10.1. Данные о росте и весе

Лицо	Рост, дюймы	Рост, сантиметры	Вес, фунты
A	63	160	150
B	67	170.2	160
C	70	177.8	171

Если измерять рост в дюймах, то ближайшим соседом *B* будет *A*:

```
a_to_b = distance([63, 150], [67, 160])      # 10.77
a_to_c = distance([63, 150], [70, 171])     # 22.14
b_to_c = distance([67, 160], [70, 171])     # 11.40
```

Однако если измерять в сантиметрах, то, напротив, ближайшим соседом *B* будет *C*:

```
a_to_b = distance([160, 150], [170.2, 160]) # 14.28
a_to_c = distance([160, 150], [177.8, 171]) # 27.53
b_to_c = distance([170.2, 160], [177.8, 171]) # 13.37
```

Очевидно, задача трудно разрешима, если смена единиц измерения приводит к таким результатам. По этой причине, в тех случаях, когда размерности друг с другом не сопоставимы, данные иногда *шкалируют*, благодаря чему каждая размерность имеет нулевое среднее значение и стандартное отклонение, равное 1. Преобразование каждой размерности в "стандартные отклонения от среднего значения" позволяет фактически избавиться от единиц измерения<sup>3</sup>.

В начале для каждого столбца требуется вычислить среднее `mean` и стандартное отклонение `standard_deviation`:

```
# стандартизировать данные
def scale(data_matrix):
    """вернуть средние и стандартные отклонения для каждого столбца"""
    num_rows, num_cols = shape(data_matrix)
    means = [mean(get_column(data_matrix, j))
              for j in range(num_cols)]
    stdevs = [standard_deviation(get_column(data_matrix, j))
              for j in range(num_cols)]
    return means, stdevs
```

<sup>3</sup> Шкалирование данных — это метод, используемый для стандартизации ряда независимых переменных или признаков и обычно выполняется на этапе предварительной обработки данных. Шкалирование также подразумевает нормирование численного значения каждой размерности так, чтобы оно попадало в интервал типа [0, 1]. — Прим. пер.

И затем использовать эту функцию для создания новой матрицы данных:

```
def rescale(data_matrix):
    """прошкалировать входящие данные так, чтобы каждый столбец
    имел нулевое среднее значение и стандартное отклонение, равное 1;
    пропускает столбцы без отклонения"""
    means, stdevs = scale(data_matrix)

    # прошкалировать по условию
    def rescaled(i, j):
        if stdevs[j] > 0:
            return (data_matrix[i][j] - means[j]) / stdevs[j]
        else:
            return data_matrix[i][j]

    num_rows, num_cols = shape(data_matrix)
    return make_matrix(num_rows, num_cols, rescaled)
```

Как всегда, нужно руководствоваться здравым смыслом. Если взять гигантский набор данных о росте и весе людей и отфильтровать его по росту между 69.5 и 70.5 дюймами, то вполне вероятно, что (в зависимости от поставленной задачи) оставшаяся вариация в данных будет представлять шум, чье стандартное отклонение, скорее всего, не стоит рассматривать на равных условиях с отклонениями других размерностей.

## Снижение размерности

Иногда "действительные" (или полезные) размерности данных могут не соответствовать имеющимся размерностям. Например, рассмотрим набор данных, изображенный на рис. 10.5.

Большая часть вариации в данных, по всей видимости, проходит вдоль единственной размерности, которая не соответствует ни оси  $x$ , ни оси  $y$ .

В таких случаях применяется один из методов факторного анализа — метод главных компонент (principal component analysis, PCA), который используется для извлечения одной или более размерностей, объясняющих большинство вариации в данных<sup>4</sup>.



На практике, этот метод не используется на наборе данных такой малой размерности. Метод снижения размерности в основном целесообразно применять к многомерному набору данных, когда требуется найти небольшое подмножество, которое объясняет большинство вариации. К сожалению, такие случаи трудно проиллюстрировать в двумерном формате книги.

<sup>4</sup> См. [http://www.machinelearning.ru/wiki/index.php?title=Метод\\_главных\\_компонент](http://www.machinelearning.ru/wiki/index.php?title=Метод_главных_компонент). — Прим. пер.

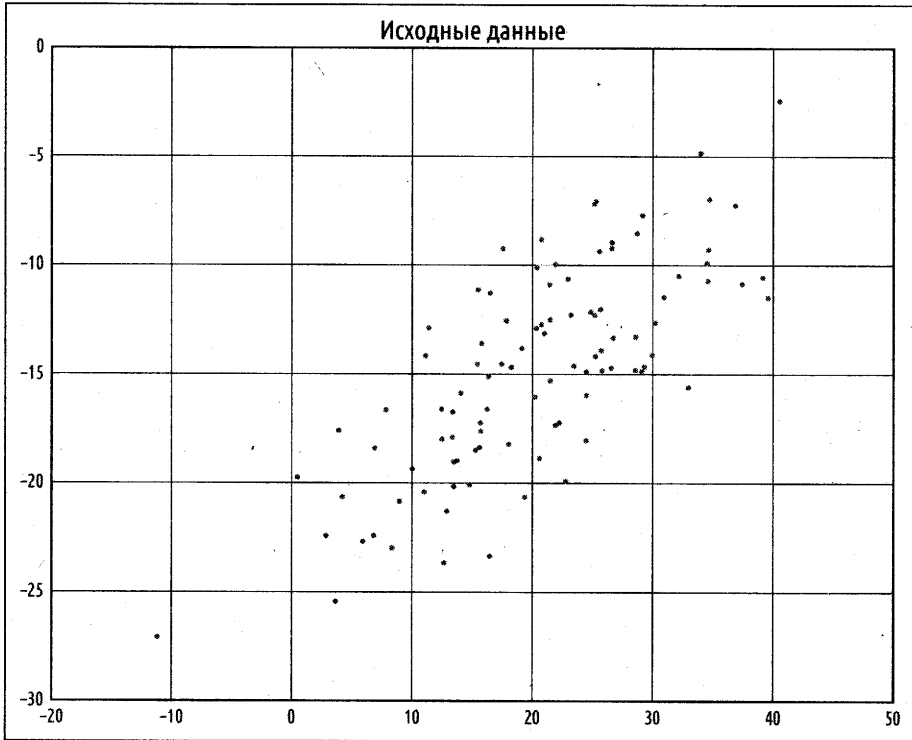


Рис. 10.5. Данные с "неправильными" осями

На первом шаге нужно пересчитать данные так, чтобы каждая размерность имела среднее нулевое (т. е. центрировать их):

```
# центрировать матрицу, т. е. вычесть из матрицы среднее значение
def de_mean_matrix(A):
    """вернуть результат вычитания из каждого значения
    в A среднего значения ее столбца. Результирующая матрица
    имеет нулевое среднее в каждом столбце"""
    nr, nc = shape(A)
    column_means, _ = scale(A)
    return make_matrix(nr, nc, lambda i, j: A[i][j] - column_means[j])
```

(Если этого не сделать, то применяемая методика, скорее всего, будет идентифицировать непосредственно среднее значение, а не вариацию в данных.)

На рис. 10.6 представлены данные из примера после вычета среднего значения.

Теперь, имея в распоряжении центрированную матрицу  $X$ , можно задаться вопросом: какое направление описывает наибольшее значение дисперсии в данных?

В частности, при заданном направлении  $\vec{d}$  (вектор единичной длины) каждая строка  $x$  в матрице простирается в размере  $\text{dot}(x, d)$  в направлении  $d$ . А каждый ненулевой вектор  $\vec{w}$  задает направление, если его прошкалировать так, чтобы он имел длину, равную 1 (т. е. нормировать):

```

# направление
def direction(w):
    mag = magnitude(w)          # длина вектора
    return [w_i / mag for w_i in w] # нормирование вектора

```

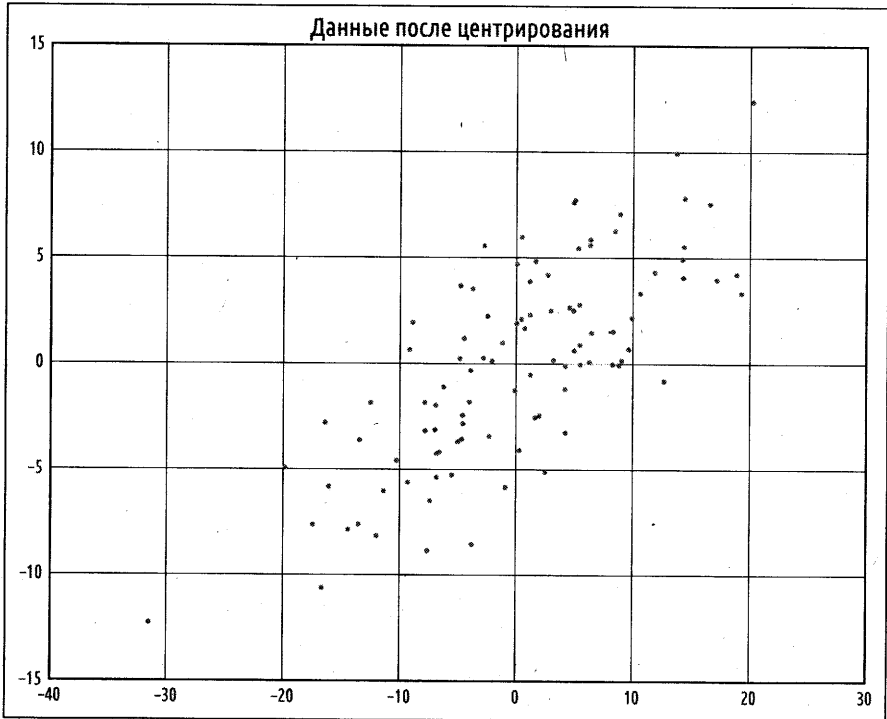


Рис. 10.6. Данные после центрирования

Следовательно, при заданном ненулевом векторе  $\vec{w}$  можно вычислить дисперсию набора данных в направлении, заданном  $w$ :

```

# направленная дисперсия строки i
def directional_variance_i(x_i, w):
    """дисперсия строки x_i в направлении, определяемом w"""
    return dot(x_i, direction(w)) ** 2

```

```

# направленная дисперсия данных
def directional_variance(X, w):
    """дисперсия данных в направлении, определяемом w"""
    return sum(directional_variance_i(x_i, w)
                for x_i in X)

```

Хотелось бы найти направление, которое максимизирует эту дисперсию. Это можно сделать методом градиентного спуска при условии, что градиент определен:

```

# градиент направленной дисперсии строки
def directional_variance_gradient_i(x_i, w):

```

```

"""вклад строки x_i в градиент w-направленной дисперсии"""
projection_length = dot(x_i, direction(w))
return [2 * projection_length * x_ij for x_ij in x_i]

# градиент направленной дисперсии данных
def directional_variance_gradient(X, w):
    return vector_sum(directional_variance_gradient_i(x_i, w)
                      for x_i in X)

```

Первая главная компонента — это как раз направление, которое максимизирует функцию направленной дисперсии `directional_variance`:

```

# первая главная компонента
def first_principal_component(X):
    guess = [1 for _ in X[0]] # предположение
    unscaled_maximizer = maximize_batch( # нешкалированный максимизатор
        partial(directional_variance, X), # теперь это функция w
        partial(directional_variance_gradient, X), # теперь это функция w
        guess)
    return direction(unscaled_maximizer)

```

Или то же самое, если использовать стохастический градиентный спуск:

```

# здесь нет оси "y", поэтому передаем вектор, состоящий из None
# и функций, которые игнорируют этот входящий аргумент
def first_principal_component_sgd(X):
    guess = [1 for _ in X[0]]
    unscaled_maximizer = maximize_stochastic(
        lambda x, _: directional_variance_i(x, w),
        lambda x, _: directional_variance_gradient_i(x, w),
        X,
        [None for _ in X], # фиктивная "y"
        guess)
    return direction(unscaled_maximizer)

```

Эти функции возвращают направление `[0.924, 0.383]` на центрированном наборе данных, которое со всей очевидностью показывает геометрическую ось, вдоль которой данные варьируют (рис. 10.7).

Найдя направление, т. е. первую главную компоненту, можно спроецировать на него данные, чтобы найти значения этой компоненты:

```

# спроецировать v на w
def project(v, w):
    """вернуть проекцию v на направление w"""
    projection_length = dot(v, w)
    return scalar_multiply(projection_length, w)

```

Если есть необходимость найти дополнительные компоненты, то из данных следует сначала удалить проекции:

```

# удалить проекцию из вектора
def remove_projection_from_vector(v, w):
    """проецирует v в w и вычитает результат из v"""
    return vector_subtract(v, project(v, w))

# удалить проекцию из данных
def remove_projection(X, w):
    """для каждой строки X
    проецирует строку в w и вычитает результат из строки"""
    return [remove_projection_from_vector(x_i, w) for x_i in X]

```

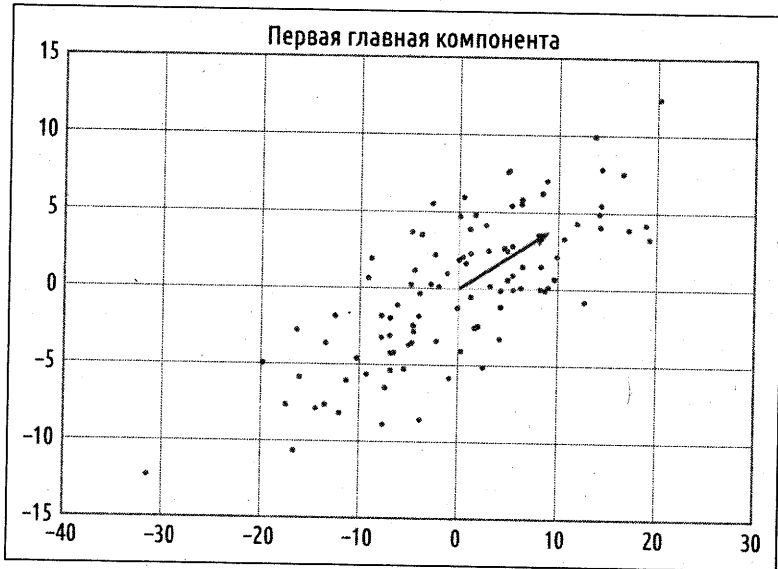


Рис. 10.7. Первая главная компонента

Поскольку в этом примере использован двумерный набор данных, то после удаления первой компоненты фактически останутся одномерные данные (рис. 10.8).

На этом этапе можно найти следующую главную компоненту, повторно применив процедуру к результату функции удаления проекции `remove_projection` (рис. 10.9).

На наборах данных большей размерности можно итеративно находить столько компонент, сколько нужно:

```

# анализ главных компонент
def principal_component_analysis(X, num_components):
    components = []
    for _ in range(num_components):
        component = first_principal_component(X)
        components.append(component)
        X = remove_projection(X, component)

    return components

```

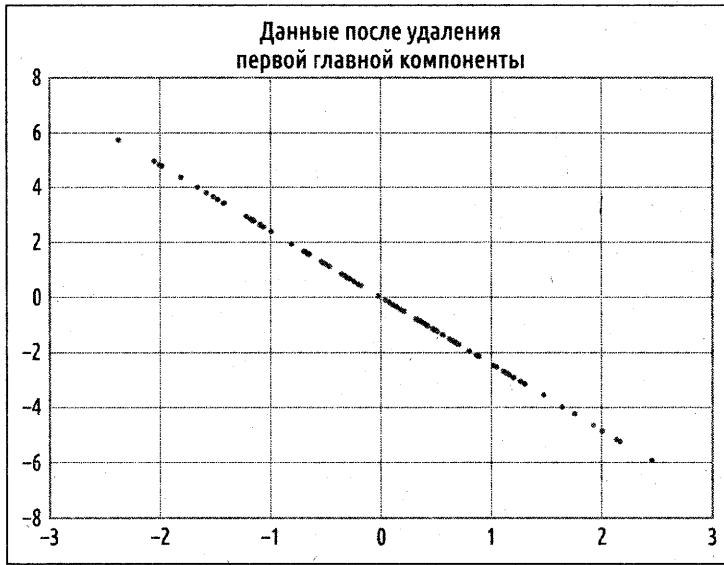


Рис. 10.8. Данные после удаления первой главной компоненты



Рис. 10.9. Первые две главные компоненты

Затем можно преобразовать данные в пространство, охватываемое компонентами, меньшей размерности:

```
# трансформировать вектор
def transform_vector(v, components):
    return [dot(v, w) for w in components]

# трансформировать табличные данные
def transform(X, components):
    return [transform_vector(x_i, components) for x_i in X]
```

Этот метод представляет ценность по нескольким причинам. Во-первых, он помогает очистить данные за счет устранения шумовых размерностей и консолидации размерностей с высокой корреляцией.

Во-вторых, после получения низкоразмерного представления данных можно использовать различные методы, которые не работают на высокоразмерных данных. Мы увидим примеры таких методов в оставшейся части книги.

Этот метод помогает усовершенствовать модели, но в то же время он усложняет их интерпретацию. Выводы наподобие того, что "каждый дополнительный год стажа увеличивает зарплату в среднем на \$10 000", абсолютно понятны, что нельзя сказать о выводах типа "любое увеличение на 0.1 единицу в третьей главной компоненте добавляет к зарплате в среднем \$10 000".

## Для дальнейшего изучения

- ◆ Как уже упоминалось в конце *главы 9*, библиотека `pandas` (<http://pandas.pydata.org/>) является, вероятно, базовым инструментом на Python для очистки, преобразования, управления и дальнейшей работы с данными. Все примеры, которые были созданы в этой главе вручную, можно реализовать гораздо проще с помощью `pandas`. Вероятно, наилучшим способом познакомиться с `pandas` является прочтение книги Уэса Маккинни "Python и анализ данных".
- ◆ Библиотека `scikit-learn` содержит целый спектр функций разложения матриц<sup>5</sup>, в том числе для проведения анализа главных компонент (PCA).

---

<sup>5</sup> <http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>.



# Машинное обучение

Всегда готов учиться, хотя не всегда нравится, когда меня учат.

Уинстон Черчилль<sup>1</sup>

Многие полагают, что наука о данных в основном имеет дело с машинным обучением и аналитики данных только и делают, что целыми днями строят, обучают и настраивают машинно-обучаемые модели. (Опять же, многие из этих людей *в действительности* не представляют, что такое машинное обучение.) По сути, наука о данных — это преимущественно сведение бизнес-проблем к проблемам в области данных, задачам сбора, понимания, очистки и форматирования данных, после чего машинное обучение идет как дополнение. Даже в этом виде оно является интересным и важным дополнением, в котором следует очень хорошо разбираться для того, чтобы решать задачи в области науки о данных.

## Моделирование

Прежде чем начать обсуждать машинное обучение, следует поговорить о *моделях*.

Что такое модель? Это, в сущности, подробное описание математической (или вероятностной) связи, которая существует между различными величинами.

Например, если необходимо собрать деньги на сайт социальной сети, то следует построить *бизнес-модель* (скорее всего, она будет в виде электронной таблицы), которая принимает входящие аргументы, такие данные, как "число пользователей", "доход от рекламы на одного пользователя" и "число сотрудников", и возвращает годовую прибыль на ближайшие несколько лет. Поваренная книга приводит к модели, которая соотносит входящие аргументы "число едоков" и "аппетит" с результирующим количеством необходимых ингредиентов. А те, кто когда-либо смотрел по телевизору покер, знают, что "вероятность победы" каждого игрока оценивается в режиме реального времени на основе модели, которая учитывает карты, раскрытые на конкретный момент времени и распределение карт в колоде.

Бизнес-модель, вероятно, будет основана на простых математических соотношениях: доходы минус расходы, где доходы — это количество проданных единиц, помноженное на среднюю цену, и т. д. Модель поваренной книги, вероятно, будет

---

<sup>1</sup> Сэр Уинстон Черчилль (1874–1965) — британский государственный и политический деятель, премьер-министр Великобритании в 1940–1945 и 1951–1955 гг.; военный (полковник), журналист, писатель. — *Прим. пер.*

основана на методе проб и ошибок — кто-то пошел на кухню и опробовал разные комбинации ингредиентов, пока не нашел то, что нравится. А модель покера будет основана на теории вероятностей, правилах игры в покер и некоторых достаточно безобидных предположениях о случайных процессах, при помощи которых карты раздаются.

## Что такое машинное обучение?

У любого специалиста имеется собственное точное определение этого термина, но мы будем использовать термин "*машинное обучение*" (machine learning), имея в виду создание и применение моделей, *извлеченных из данных в результате обучения*. В других контекстах это может называться *прогноznым моделированием* (predictive modeling) и *интеллектуальным анализом данных* (data mining), но мы будем придерживаться термина "машинное обучение". Как правило, задача будет заключаться в том, чтобы использовать существующие данные для разработки моделей, которые можно применять для *предсказания* разных исходов в отношении новых данных, как например, предсказание:

- ◆ является ли сообщение спамом или нет;
- ◆ является ли транзакция по кредитной карте мошеннической;
- ◆ на каких рекламных сообщениях покупатель с наибольшей вероятностью будет нажимать;
- ◆ какая футбольная команда выиграет Суперкубок.

Мы обратимся к моделям с *учителем* (т. е. таким, где есть выборка, маркированная правильными ответами, на которых проходит обучение) и моделям *без учителя* (где правильные ответы не заданы). Кроме них существуют и другие типы моделей, такие как модели с *частичным* обучением (где только некоторые данные маркированы правильными ответами) и модели с *динамическим* обучением (где модель должна постоянно приспосабливаться ко вновь прибывающим данным), которые мы не будем рассматривать в этой книге.

При этом, даже в самой простой ситуации интересующие нас связи могут быть описаны целостными совокупностями или семействами *параметризованных* моделей, и в большинстве случаев мы сами выбираем такое семейство моделей, а затем используем данные, чтобы вычислить параметры, которые в определенной мере являются оптимальными.

Например, можно предположить, что рост человека является (приблизительно) линейной функцией его веса, и затем использовать данные, чтобы вычислить вид этой линейной функции. Либо предположить, что дерево принятия решений является хорошим механизмом для диагностики заболеваний пациентов, и потом использовать данные, чтобы вычислить подобное "оптимальное" дерево. В остальной части книги мы займемся изучением различных семейств моделей, которые можно вычислить.

Но прежде чем этим заняться, следует получше разобраться в основах машинного обучения. В оставшейся части главы мы обсудим некоторые из этих базовых понятий, прежде чем перейдем к самим моделям.

## Переобучение и недообучение

Общеизвестной ловушкой в машинном обучении является *переобучение* (overfitting), приводящее к прогнозной модели, которая хорошо работает на данных из обучающей выборки, но при этом плохо обобщает на любых новых данных.

Это явление может быть вызвано тем, что модель обучена *шуму* в данных либо обучена различать конкретные входящие значения, а не те факторы, которые на самом деле несут для требуемого результата предсказательный характер.

Другую сторону этого явления представляет *недообучение* (underfitting), которое приводит к прогнозной модели, не работающей даже на обучающей выборке. Впрочем, когда это случается, обычно делают вывод, что модель недостаточно хороша, и продолжают искать более подходящую.

На рис. 11.1 выполнена полиномиальная аппроксимация выборки из данных. (Не стоит сейчас озадачиваться тем, как это сделать; мы вернемся к этому в последующих главах.)

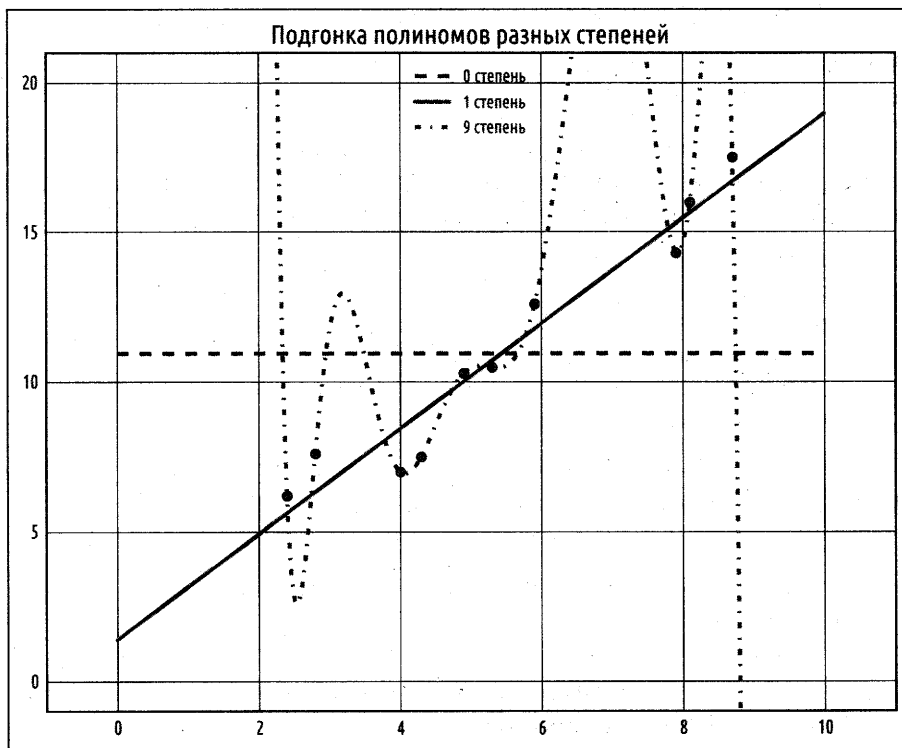


Рис. 11.1. Переобучение и недообучение

Горизонтальная прямая соответствует наиболее подходящему значению полинома 0-й степени (т. е. постоянного). Она крайне *недостаточно подогнана* под обучающую выборку (т. е. *недообучена*). Наиболее подходящее значение кривой полинома 9-й степени (т. е. из 10 параметров) проходит ровно через каждую точку обучающей выборки, но оно *чрезмерно подогнано* под выборку (т. е. *переобучено*) — если бы пришлось добавить еще несколько точек, то кривая, вполне вероятно, прошла бы мимо на большом от них удалении. А вот линия полинома 1-й степени (линейная интерполяция) достигает хорошего равновесия — она находится довольно близко к каждой точке, и (если данные репрезентативны, то) прямая, скорее всего, будет находиться близко и к новым точкам данных.

Ясно, что слишком сложные модели приводят к переобучению и плохо обобщают за пределами данных, на которых они обучены. Тогда как сделать так, чтобы модели не были слишком сложными? Наиболее фундаментальный подход предполагает использование разных данных для обучения модели и для проверки ее качества.

Простейший способ состоит в разбиении набора данных, благодаря чему (например) две трети используются для обучения модели, после чего производительность модели измеряется на оставшейся трети:

```
# сегментировать данные
def split_data(data, prob):
    """разбиение данных на части [prob, 1 - prob]"""
    results = [], []
    for row in data:
        results[0 if random.random() < prob else 1].append(row)
    return results
```

Часто в нашем распоряжении будут матрица  $x$ , состоящая из входящих переменных, и вектор  $y$  — из выходящих переменных. В этом случае необходимо разместить соответствующие значения либо в обучающей, либо в контрольной выборке:

```
# разбиение на обучающую и контрольную выборки
def train_test_split(x, y, test_pct):
    data = zip(x, y) # объединить соответствующие значения
    train, test = split_data(data, 1 - test_pct) # разбить список пар
    x_train, y_train = zip(*train) # трюк с разъединением списков
    x_test, y_test = zip(*test)
    return x_train, x_test, y_train, y_test
```

Благодаря этому можно сделать что-то наподобие следующего:

```
model = SomeKindOfModel() # некая модель
x_train, x_test, y_train, y_test = train_test_split(xs, ys, 0.33) # разбить
model.train(x_train, y_train) # обучить
performance = model.test(x_test, y_test) # проверить производительность
```

Если модель переобучена, то можно надеяться, что она покажет очень слабые результаты на (совершенно автономных) контрольных данных. Говоря иначе, если она хорошо работает на контрольных данных, то уровень уверенности в том, что модель *обучена* нежелательно *переобучена*, гораздо выше.

Тем не менее, в нескольких случаях это может оказаться неверным.

Во-первых, в контрольной и обучающей выборках могут иметься типичные схемы (паттерны), которые не обобщаются на более крупный набор данных.

Например, пусть набор данных состоит из действий пользователей, по одной строке на каждого пользователя в неделю. В таком случае большинство пользователей будут появляться как в обучающей, так и в контрольной выборках, и, как результат, некоторые модели научатся *определять* пользователей, а не искать связи с участием *атрибутов*. Данная проблема не приводит к большим неприятностям, хотя однажды мне пришлось с ней столкнуться.

Во-вторых, и это представляет более значительную проблему, когда разбивка на контрольную и обучающую выборки используется не только для того, чтобы судить о модели, но и для того, чтобы *выбирать* среди нескольких моделей. В этом случае, несмотря на то, что каждая отдельная модель может оказаться непереобученной, обучающее метаправило, которое заставляет "выбирать ту модель, которая на контрольной выборке демонстрирует наилучшие результаты", превращает функцию контрольной выборки во вторую обучающую выборку. (И естественно, модель, которая на контрольной выборке продемонстрировала наилучшие результаты, продолжит демонстрировать на контрольной выборке хорошие результаты.)

В такой ситуации необходимо разделить данные на три части: *обучающее* подмножество для построения моделей, *валидационное* подмножество (подмножество перекрестной проверки) для выбора одной из обученных моделей и *контрольное* подмножество для оценивания окончательной модели.

## Правильность модели

Когда я не занимаюсь аналитической работой, то люблю копаться в медицине. В свободное время я придумал дешевый и неинвазивный тест, который можно проводить у новорожденных. Он предсказывает предрасположенность новорожденного к развитию лейкоза с точностью свыше 98%. Мой адвокат убедил меня, что тест не патентоспособный, поэтому поделюсь подробностями: модель предсказывает лейкоз, если и только если ребенка зовут Люк (звучит похоже).

Как мы увидим ниже, этот тест действительно более чем на 98% точен. И тем не менее, он невероятно дурацкий, являясь хорошей иллюстрацией к тому, почему обычно критерий "правильности" (correctness) не используют для измерения качества модели.

Пусть строится модель, которая должна формулировать *бинарные* суждения. Является ли это письмо спамом? Следует ли нанять этого претендента? Является ли этот авиапассажир скрытым террористом?

При наличии набора маркированных данных и такой прогнозной модели каждая точка данных принадлежит одной из четырех категорий:

- ♦ истинноположительная: "это сообщение спамное, и его спамность была предсказана правильно";

- ◆ ложноположительная (ошибка 1-го рода): "это сообщение не спамное, однако была предсказана его спамность";
- ◆ ложноотрицательная (ошибка 2-го рода): "это сообщение спамное, однако была предсказана его неспамность";
- ◆ истинноотрицательная: "это сообщение не спамное, и его неспамность была предсказана правильно".

Эти категории нередко представляют в виде показателей в *таблице сопряженности* (confusion matrix) — табл. 11.1.

Таблица 11.1. Таблица сопряженности

Предположение	Спам	Не спам
Предсказано "спам"	Истинноположительная	Ложноположительная
Предсказано "не спам"	Ложноотрицательная	Истинноотрицательная

Посмотрим, как тест на лейкоз вписывается в эти рамки. В наше время приблизительно 5% младенцам из 1000 дают имя Люк<sup>2</sup>, а распространенность лейкоза на протяжении жизни составляет около 1,4%, или 14 человек на каждую 1000.

Если, будучи уверенным, что эти два фактора являются независимыми, применить мой так называемый тест "Люк означает лейкоз" к 1 млн человек, то можно ожидать, что таблица сопряженности будет выглядеть так, как табл. 11.2.

Таблица 11.2. Таблица сопряженности для теста

Предположение	Лейкоз	Не лейкоз	Всего
"Люк"	70	4930	5000
"Не Люк"	13 930	981 070	995 000
Всего	14 000	986 000	1 000 000

Затем этими данными можно воспользоваться для вычисления разного рода статистик, касающихся производительности модели. Например, вот показатель *точности* (accuracy), который определен в виде доли правильных предсказаний:

# точность

```
def accuracy(tp, fp, fn, tn):
    correct = tp + tn # правильный, если истинноположит. и истинноотрицат.
    total = tp + fp + fn + tn
    return correct / total
```

```
print(accuracy(70, 4930, 13930, 981070)) # 0.98114
```

<sup>2</sup> <http://www.babycenter.com/babyNameAllPops.htm?babyNameId=2918>.

И, как можно убедиться, он показывает впечатляющие результаты. Однако с самого начала было совершенно очевидно, что этот тест неправильный, а значит, вероятно, совсем не стоит доверять такому грубому показателю точности.

В соответствии со стандартной практикой принято обращаться к сочетанию *точности* (или прецизионности, *precision*) и *полноты* (*recall*). При этом здесь точность — это доля *истинноположительных* предсказаний относительно всех положительных предсказаний:

```
# точность (как степень положительных значений прогнозов)
def precision(tp, fp, fn, tn):
    return tp / (tp + fp)
```

```
print(precision(70, 4930, 13930, 981070)) # 0.014
```

А полнота — это доля положительных предсказаний, которую модель идентифицировала:

```
# полнота
def recall(tp, fp, fn, tn):
    return tp / (tp + fn)
```

```
print(recall(70, 4930, 13930, 981070)) # 0.005
```

Оба результата ужасны и отражают тот факт, что сама модель ужасна.

В некоторых случаях точность и полнота объединяются в метрику *F1*, которая определена следующим образом:

```
# метрика F1
def f1_score(tp, fp, fn, tn):
    p = precision(tp, fp, fn, tn)
    r = recall(tp, fp, fn, tn)

    return 2 * p * r / (p + r)
```

Она представляет собой *гармоническое среднее значение*<sup>3</sup> точности и полноты и с неизбежностью лежит между ними.

В большинстве же случаев выбор модели зависит от компромисса между точностью и полнотой. Модель, которая предсказывает "да" с небольшим уровнем уверенности, наверное, будет иметь высокую полноту и низкую точность, тогда как модель, которая дает такое же предсказание только с крайне высоким уровнем уверенности, вероятнее всего, будет иметь низкую полноту и высокую точность.

С другой стороны, этот компромисс можно представить, как некий баланс между ложноположительными и ложноотрицательными категориями. Ответ "да" слишком часто будет давать много ложноположительных результатов, а ответ "нет" — много ложноотрицательных.

<sup>3</sup> [https://en.wikipedia.org/wiki/Harmonic\\_mean](https://en.wikipedia.org/wiki/Harmonic_mean), [https://ru.wikipedia.org/wiki/Среднее\\_гармоническое](https://ru.wikipedia.org/wiki/Среднее_гармоническое).

Предположим, что имеется 10 факторов риска лейкемии, и чем больше их у человека идентифицировано, тем больше шансов, что должен развиться лейкоз. В таком случае будет существовать непрерывное пространство тестов: "предсказать лейкоз при не менее одном факторе", "предсказать лейкоз при не менее двух факторах" и т. д. Если увеличивать пороговую величину, то точность теста будет расти (поскольку люди с более высоким фактором риска более склонны к развитию заболевания), а полнота теста будет уменьшаться (т. к. все меньше потенциальных больных достигнут порога). В таких случаях выбор правильной пороговой величины становится вопросом нахождения правильного компромисса.

## Компромисс между смещением и дисперсией

Еще один способ взглянуть на проблему переобучения заключается в нахождении баланса между смещением и дисперсией.

Оба показателя позволяют измерить результаты многократного переучивания модели на разных обучающих выборках (из той же самой более крупной или генеральной совокупности).

Например, модель с полиномом 0-й степени из *разд. "Переобучение и недообучение"* данной главы сделает много ошибок практически на любой обучающей выборке (взятой из той же самой совокупности), а значит, имеет высокое *смещение*. Однако любые две обучающие выборки, взятые случайно, дадут вполне похожие модели (поскольку имеют вполне похожие средние значения). И поэтому говорят, что модель имеет низкую *дисперсию*. А высокое смещение и низкая дисперсия, как правило, соответствуют недообучению.

С другой стороны, модель с полиномом 9-й степени, которая хорошо подогнана под обучающую выборку, имеет очень низкое смещение, но очень высокую дисперсию (поскольку любые две обучающие выборки, скорее всего, породят очень разные модели). А это соответствует переобучению.

Обдумывание модельных задач, таким образом, помогает понять, что делать, когда модель не работает, как надо.

Если модель имеет высокое смещение (и значит, дает плохие результаты даже на обучающих выборках), то следует попробовать *добавить* больше признаков. Переход от модели с полиномом 0-й степени к модели с полиномом 1-й степени в результате дал большое улучшение.

Если модель имеет высокую дисперсию, то можно точно также *удалить* признаки. Впрочем, еще одно решение состоит в том, чтобы добыть больше данных (если это возможно).

На рис. 11.2 подобран полином 9-й степени к выборкам различных объемов. На основе 10 точек данных модель подогнана по всему полю, как было показано ранее. Если, напротив, обучить модель на 100 точках данных, то переобученность noticeably снизится. А модель, обученная на основе 1000 точек, будет выглядеть очень похожей на модель на основе линейной интерполяции (с полиномом 1-й степени).



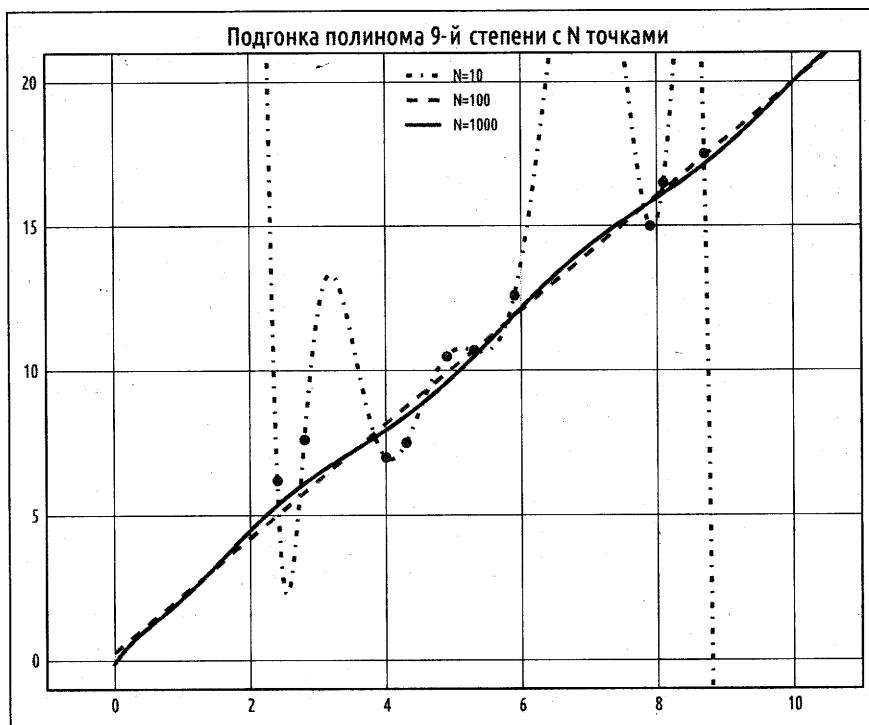


Рис. 11.2. Уменьшение дисперсии при помощи дополнительных данных

Если держать сложность модели постоянной, то чем больше данных, тем труднее ее переобучить.

С другой стороны, дополнительные данные не улучшат смещение. Если модель не использует достаточного количества признаков для объяснения закономерностей в данных, то мобилизация дополнительных данных не поможет.

## Извлечение и отбор признаков

Как уже упоминалось, когда данные не содержат достаточного количества признаков, то модель, скорее всего, будет недообучена. Когда же признаков слишком много, то ее легко переобучить. Но что такое признаки и откуда они берутся?

*Признаки* — это любые входящие значения, которые передаются в модель.

В простейшем случае признаки просто заданы. Если нужно предсказать чью-либо зарплату на основе его многолетнего стажа, то число лет будет единственным предполагаемым признаком.

(Впрочем, как мы видели в разд. "Переобучение и недообучение" данной главы, помимо этого можно попробовать добавить стаж в квадрате, в кубе и т. д. в случае, если это помогает улучшить модель.)

По мере усложнения данных ситуация становится интереснее. Представим, что строится спам-фильтр, который должен предсказывать, является ли сообщение

электронной почты нежелательным или нет. Большинство моделей не готовы работать с необработанным почтовым сообщением, которое представляет собой всего лишь фрагмент текста. Вам следует извлечь признаки. Например:

- ◆ Присутствует ли в почтовом сообщении слово "виагра"?
- ◆ Сколько раз появляется буква d?
- ◆ Каков домен отправителя?

Первый признак будет иметь лишь два значения — "да" или "нет", которые обычно кодируются как 1 или 0; второй — это число, а третий — альтернатива из дискретного набора вариантов.

Извлекаемые из данных признаки почти всегда попадают в одну из этих трех категорий. Более того, тип имеющихся признаков накладывает ограничения на тип моделей, которые можно использовать.

Наивный байесовский классификатор, который будет построен в *главе 13*, подходит для бинарных признаков (типа да/нет), как и первый в предыдущем списке.

Для регрессионных моделей, которые будут рассмотрены в *главах 14 и 16*, требуются численные признаки (которые могут включать в себя фиктивные переменные, представляющие собой нули и единицы).

Модели деревьев принятия решений, к которым мы обратимся в *главе 17*, могут иметь дело как с численными, так и с категориальными данными.

В примере со спам-фильтром были рассмотрены приемы создания признаков. Вместе с тем иногда, напротив, требуется найти способы их удаления.

Например, входящими значениями могут являться векторы из нескольких сотен чисел. В зависимости от ситуации, возможно, будет целесообразней свести их к нескольким важным размерностям (как в *разд. "Снижение размерности" главы 10*) и использовать только это небольшое число признаков. Или же использовать метод (например, регуляризацию, к которой мы обратимся в *разд. "Регуляризация" главы 15*), который штрафует модели тем больше, чем больше признаков в них используется.

Как выполняется отбор признаков? Здесь в игру вступает сочетание *опыта* и *знания* предметной области. Если получено много электронных писем, то вполне вероятно, что появится некое представление о том, что присутствие некоторых слов может быть хорошим индикатором спамности, а также что число вхождений буквы d в текст, скорее всего, не является его хорошим индикатором. Но в целом необходимо опробовать разные подходы, что является частью игры.

## Для дальнейшего изучения

- ◆ Продолжайте читать! Следующие несколько глав посвящены различным семьям моделей машинного обучения.
- ◆ Курс по машинному обучению (<https://www.coursera.org/learn/machine-learn>) в рамках онлайн-платформы Courser Стэнфордского университета является

своеобразным массовым открытым онлайн-курсом (МООК) и хорошей площадкой для более глубокого постижения основ машинного обучения. МООК по машинному обучению (<https://work.caltech.edu/telecourse.html>) в Калифорнийском технологическом институте тоже подойдет.

- ◆ Учебник "Элементы статистического обучения" является несколько каноническим пособием, которое можно скачать бесплатно (<http://statweb.stanford.edu/~tibs/ElemStatLearn/>). Но будьте осторожны: там *много* математических выкладок.

## К ближайших соседей

Если желаешь досадить соседям, расскажи им правду о них.

Пьетро Аретино<sup>1</sup>

Представим, что поставлена задача предсказать, как я буду голосовать на следующих президентских выборах. Если обо мне (кроме данных о месте проживания) больше ничего не известно, то целесообразно посмотреть на то, как собираются голосовать мои *соседи*. Живя, как и я, в центре Сиэтла, они неизменно планируют голосовать за демократического кандидата, что позволяет с большой долей уверенности предположить, что я тоже сделаю свой выбор в пользу кандидата от демократов.

Теперь, предположим, обо мне известно больше, чем просто география — возможно, известен мой возраст, доход, сколько у меня детей и т. д. В той мере, в какой мое поведение обуславливается (или характеризуется) этими признаками, рассмотрение только тех соседей, которые расположены ко мне близко среди всех этих размерностей, скорее всего, позволит сделать более точное предсказание, чем рассматривать абсолютно всех соседей. В этом и заключается основная идея метода *классификации на основе ближайших соседей* (nearest neighbors classification).

### Модель

Прогнозная модель классификации на основе *ближайших соседей* является одной из простейших. В ней не делается никаких математических допущений и не требуется какого-то тяжелого функционального аппарата. Единственное, что требуется, это:

- ◆ некоторое представление о расстоянии;
- ◆ предположение, что точки, расположенные друг к другу близко, подобны.

Большинство методов, которые рассматриваются в этой книге, обращаются ко всей выборке в целом с тем, чтобы в результате обучения извлечь из данных схемы или закономерности. В отличие от них метод ближайших соседей вполне осознанно пренебрегает значительной частью информации, поскольку предсказание для любой новой точки зависит всего лишь от нескольких ближайших к ней точек.

<sup>1</sup> Пьетро Аретино (1492–1556) — итальянский писатель Позднего Ренессанса, сатирик, публицист и драматург, считающийся некоторыми исследователями предтечей и основателем европейской журналистики (см. [https://ru.wikipedia.org/wiki/Пьетро\\_Аретино](https://ru.wikipedia.org/wiki/Пьетро_Аретино)). — *Прим. пер.*

Более того, ближайшие соседи, скорее всего, не помогут понять движущие силы любого изучаемого явления. Попытка предсказать мое голосование на выборах на основе голосования моих соседей мало что скажет о том, что заставляет меня голосовать именно таким образом, тогда как некая альтернативная модель, которая предсказывает голосование на основе (допустим) дохода и семейного положения, очень даже может.

В общей ситуации, имеются некие точки данных и соответствующий набор меток. Метки могут быть True и False, обозначая, удовлетворяет ли каждое входящее значение некоторому условию, такому как "спам?" или "ядовитый?" или "приятный для просмотра?" Либо они могут быть категориями, такими как рейтинги фильмов (G, PG, PG-13, R, NC-17), либо именами кандидатов в президенты, либо названиями предпочтительных языков программирования.

В нашем случае точки данных будут представлены векторами, вследствие чего можно воспользоваться функцией расстояния `distance` из главы 4.

Пусть число  $k = 3$  или 5. Тогда, если необходимо классифицировать некоторые новые точки данных, то следует найти  $k$  ближайших маркированных точек и дать им проголосовать за новый результат.

Для этого потребуется функция, которая подсчитывает голоса. Одна из возможностей заключается в следующем:

```
# грубый подсчет большинства голосов
def raw_majority_vote(labels):
    votes = Counter(labels)
    winner, _ = votes.most_common(1)[0]
    return winner
```

Однако она ничего не предпринимает в ситуации равного числа голосов. Например, допустим, фильмам присваиваются рейтинги, и пяти ближайшим фильмам присвоены рейтинги G, G, PG, PG и R, где фильмы с рейтингами G и PG имеют по два голоса. В таком случае, имеется несколько вариантов:

- ◆ выбрать одного из победителей случайным образом;
- ◆ взвесить голоса по удаленности и выбрать взвешенного победителя;
- ◆ уменьшать  $k$ , пока не будет найден единственный победитель.

Реализуем третий вариант:

```
# отбор по большинству голосов
def majority_vote(labels):
    """подразумевает, что метки упорядочены от ближайшей
    до самой удаленной"""
    vote_counts = Counter(labels)
    winner, winner_count = vote_counts.most_common(1)[0]
    num_winners = len([count
                       for count in vote_counts.values()
                       if count == winner_count])
```

```

if num_winners == 1:
    return winner # единственный победитель, поэтому вернуть его
else:
    return majority_vote(labels[: -1]) # попытаться снова
                                        # без самого дальнего

```

Такой подход обязательно в конце концов сработает, т. к. в худшем случае все сведется только к одной метке, которая и будет победителем.

При помощи этой функции можно легко создать классификатор:

```

# классифицировать на основе k ближайших соседей
def knn_classify(k, labeled_points, new_point):
    """каждая маркированная точка должна быть представлена
    парой (точка, метка)"""

    # упорядочить маркированные точки от ближайшей до самой удаленной
    by_distance = sorted(labeled_points,
                        key=lambda (point, _): distance(point, new_point))

    # найти метки для k ближайших
    k_nearest_labels = [label for _, label in by_distance[:k]]

    # и дать им проголосовать
    return majority_vote(k_nearest_labels)

```

Посмотрим, как это работает.

## Пример: предпочтительные языки

Результаты первого опроса пользователей социальной сети DataSciencester готовы. Они содержат данные о языках программирования, наиболее предпочитаемых среди пользователей сети в ряде крупных городов:

```

# города
# каждая запись - это ([долгота, широта], предпочтительный язык)
cities = [[(-122.3, 47.53), "Python"], # Сиэтл
          [(-96.85, 32.85), "Java"],   # Остин
          [(-89.33, 43.13), "R"],      # Мэдисон
          # ... и т. д.
]

```

Директор по вовлечению общественности хотел бы знать, можно ли применить эти результаты для того, чтобы предсказать наиболее предпочитаемые языки программирования для тех мест, которые не были частью исследования.

Как всегда, первым делом надо вывести данные на диаграмму (рис. 12.1):

```

# ключ = язык, значение = пара (долгота, широта)
plots = { "Java" : ([], []), "Python" : ([], []), "R" : ([], []) }

```

```

# каждый язык должен иметь разную метку и цвет
markers = { "Java" : "o", "Python" : "s", "R" : "^" }
colors = { "Java" : "r", "Python" : "b", "R" : "g" }

for (longitude, latitude), language in cities:
    plots[language][0].append(longitude)
    plots[language][1].append(latitude)

# создать серии разброса для каждого языка
for language, (x, y) in plots.iteritems():
    plt.scatter(x, y, color=colors[language], marker=markers[language],
               label=language, zorder=10)

plot_state_borders(plt) # отрисовать границы штатов
                        # претворимся, что есть функция, которая это делает
plt.legend(loc=0) # пусть matplotlib сама выберет расположение легенды
plt.axis([-130,-60,20,55]) # установить оси

plt.title("Наиболее предпочитаемые языки")
plt.show()

```



Рис. 12.1. Наиболее предпочитаемые языки программирования



Обращаем внимание на вызов функции `plot_state_borders()`, которая на самом деле еще не определена. Несмотря на то, что ее реализация имеется на странице GitHub<sup>2</sup> данной книги, все же стоит в качестве упражнения попробовать реализовать ее самостоятельно.

<sup>2</sup> <https://github.com/joelgrus/data-science-from-scratch>.

Для этого надо:

1. Найти в Интернете соответствующие данные по запросу *широта долгота государственной границы*.
2. Конвертировать любые данные, которые будут найдены, в список отрезков [(долг1, шир1), (долг2, шир2)].
3. Использовать `plt.plot()` для отрисовки отрезков.

Поскольку, по всей видимости, в близлежащих местах, как правило, предпочитают один и тот же язык, то метод  $k$  ближайших соседей представляется разумным вариантом прогнозной модели.

Для начала посмотрим, что произойдет, если попытаться предсказать наиболее предпочитаемый язык для каждого города с помощью его соседей, исключая рассматриваемый город:

```
# попробовать несколько разных значений для k
for k in [1, 3, 5, 7]:
    num_correct = 0

    for city in cities:
        location, actual_language = city
        other_cities = [other_city
                        for other_city in cities
                        if other_city != city]

        # предсказанный язык
        predicted_language = knn_classify(k, other_cities, location)

        if predicted_language == actual_language: # равен фактическому
                                                    # языку?
            num_correct += 1

    print(k, "сосед(а,ей):", num_correct, "правильных из", len(cities))
```

Похоже, алгоритм работает лучше всего при  $k=3$ , давая правильный результат примерно в 59% случаев:

```
Количество соседей 1: 40 правильных из 75
Количество соседей 3: 44 правильных из 75
Количество соседей 5: 41 правильных из 75
Количество соседей 7: 35 правильных из 75
```

Теперь можно взглянуть на характер распределения языков по регионам в условиях каждой схемы ближайших соседей. Это делается путем классифицирования всей сетки точек, после чего строится диаграмма, так же как в случае с городами:

```
plots = { "Java" : ([], []), "Python" : ([], []), "R" : ([], []) }

k = 1 # или 3 или 5 или ...

for longitude in range(-130, -60):
    for latitude in range(20, 55):
        predicted_language = knn_classify(k, cities, [longitude, latitude])
```



```
plots[predicted_language][0].append(longitude)
plots[predicted_language][1].append(latitude)
```

К примеру, на рис. 12.2 показано, что происходит, если обратиться лишь к одному ближайшему соседу ( $k=1$ ).



Рис. 12.2. Языки программирования на основе одного ближайшего соседа



Рис. 12.3. Языки программирования на основе трех ближайших соседей

Видно много резких переходов от одного языка к другому с ломаными границами. При увеличении числа соседей до 3 границы регионов становятся более гладкими для каждого языка (рис. 12.3).

При увеличении числа соседей до 5 границы становятся еще более гладкими (рис. 12.4).

В рассматриваемом примере размерности примерно сопоставимы. Однако если это не так, то их следует прошкалировать так же, как в разд. "Многомерное шкалирование" главы 10.

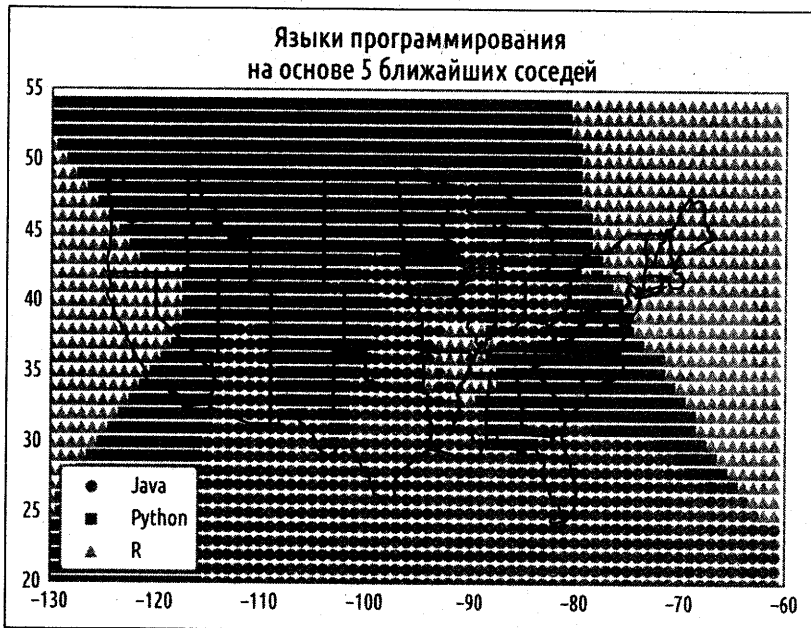


Рис. 12.4. Языки программирования на основе пяти ближайших соседей

## Проблема проклятия размерности

В более высоких размерностях метод  $k$  ближайших соседей сталкивается с трудностями из-за проблемы "проклятия размерности", которая сводится к тому, что высокоразмерные пространства обширны. Точки в таких пространствах, как правило, вообще не располагаются близко друг к другу. Чтобы убедиться в этом, надо сгенерировать пары случайных точек в разных размерностях  $d$ -мерного "единичного куба" и вычислить между ними расстояния.

Генерирование случайных точек должно уже войти в привычку:

```
# генерирование случайных точек
def random_point(dim):
    return [random.random() for _ in range(dim)]
```

так же как и написание функции, которая генерирует расстояния:

```
# генерирование случайных расстояний
def random_distances(dim, num_pairs):
    return [distance(random_point(dim), random_point(dim))
            for _ in range(num_pairs)]
```

Для каждой размерности от 1 до 100 вычислим 10 000 расстояний и воспользуемся ими для того, чтобы вычислить между точками среднее и минимальное расстояния в каждой размерности (рис. 12.5):

```
dimensions = range(1, 101) # размерности

avg_distances = [] # средние расстояния
min_distances = [] # минимальные расстояния

random.seed(0)
for dim in dimensions:
    distances = random_distances(dim, 10000) # 10 000 произвольных пар
    avg_distances.append(mean(distances)) # отследить средние
    min_distances.append(min(distances)) # отследить минимальные
```

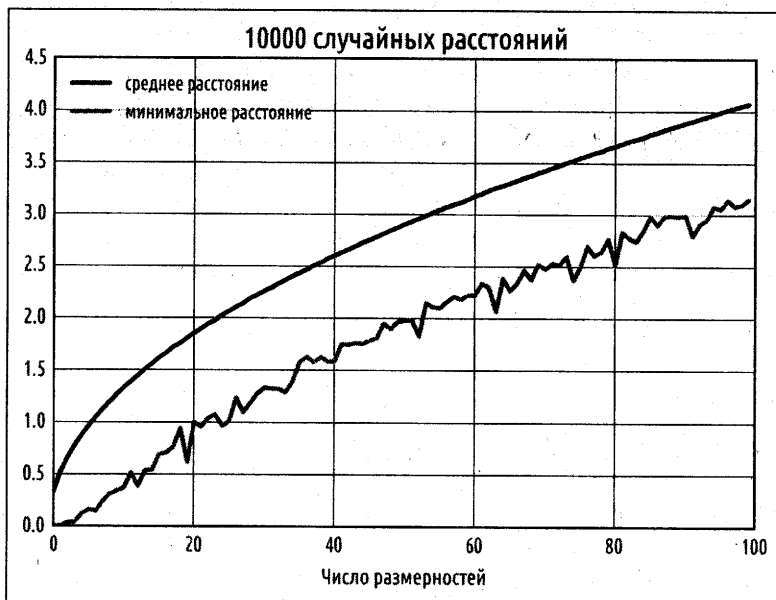


Рис. 12.5. Проблема проклятия размерности

При увеличении размерности среднее расстояние между точками увеличивается. Но больше проблем вызывает соотношение между минимальным и средним расстояниями (рис. 12.6):

```
min_avg_ratio = [min_dist / avg_dist
                 for min_dist, avg_dist in zip(min_distances, avg_distances)]
```



Рис. 12.6. Проблема проклятия размерности (повторно)

В низкоразмерных наборах данных ближайшие точки, как правило, расположены гораздо ближе средних. Однако две точки располагаются близко, только если они расположены близко в каждой размерности, а каждая дополнительная размерность — даже если это всего лишь шум — предоставляет еще одну возможность для того, чтобы любая точка стала еще дальше от любой другой точки. При большом числе размерностей ближайшие точки могут быть не намного ближе, чем средние, вследствие чего близкая расположенность двух точек теряет особый смысл (если только данные не структурированы, что может заставить их вести себя так, как если бы они обладали значительно меньшей размерностью).

Другой способ взглянуть на эту проблему касается свойства разреженности пространств более высокой размерности.

Если выбрать 50 случайных чисел между 0 и 1, то скорее всего получится довольно хорошая выборка на основе единичного интервала (рис. 12.7).

Если выбрать 50 случайных точек в единичном квадрате, то покрытие выборки будет меньше (рис. 12.8).

И еще меньше в трех измерениях (рис. 12.9).

Библиотека `matplotlib` не очень хорошо справляется с диаграммами в четырех измерениях, так что остановимся на трех, но уже сейчас можно заметить, что начинают появляться большие пустые пространства без точек возле них. При увеличении размерности — если только данные не растут в геометрической прогрессии — эти большие пустоты представляют собой области, удаленные от всех тех точек, которые нужны для прогнозирования.

В целом, если метод ближайших соседей применяется к данным более высокой размерности, то, вероятно, сперва стоит каким-то образом снизить их размерность.

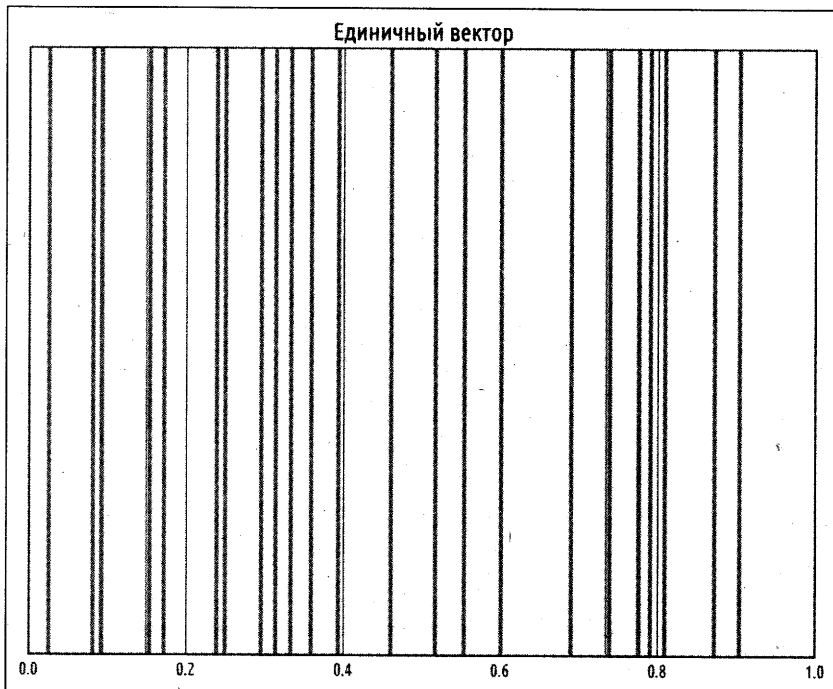


Рис. 12.7. Пятьдесят случайных точек в одном измерении

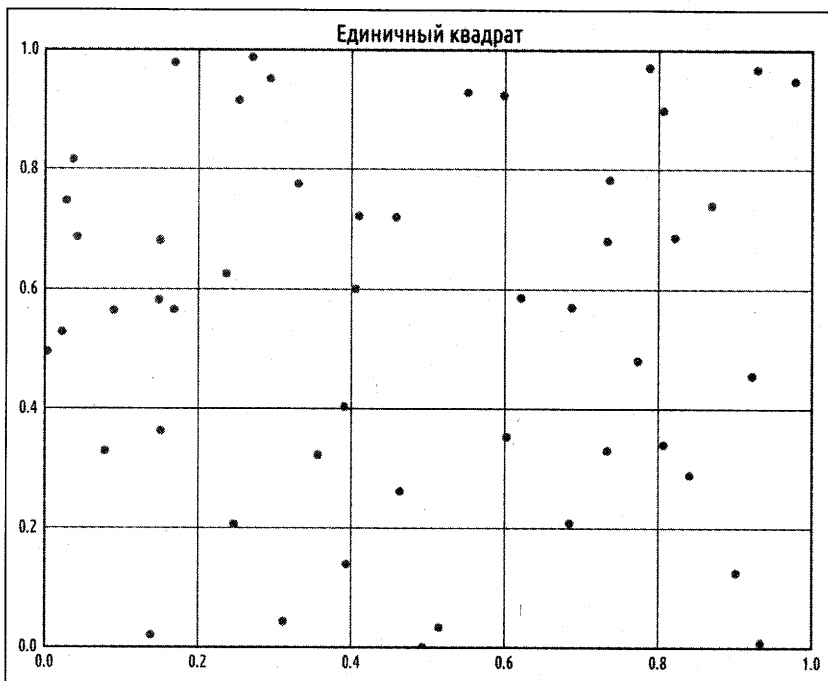


Рис. 12.8. Пятьдесят случайных точек в двух измерениях

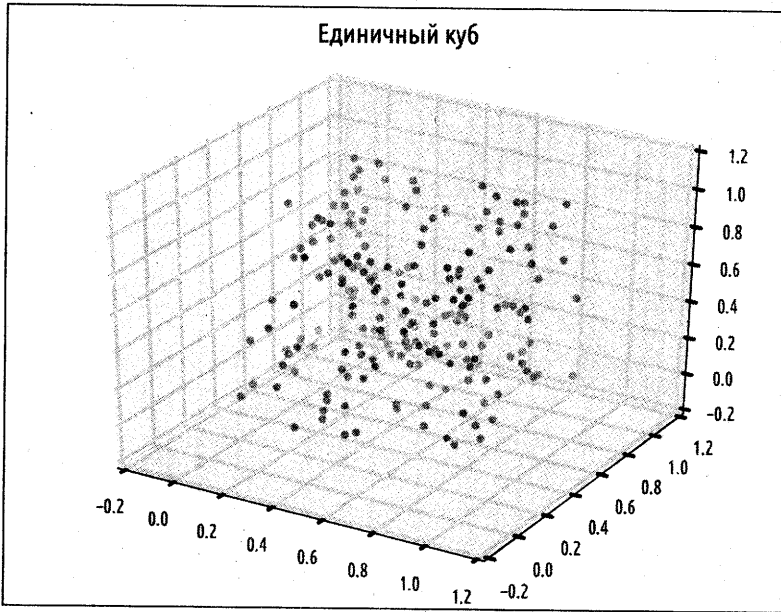


Рис. 12.9. Пятьдесят случайных точек в трех измерениях

## Для дальнейшего изучения

Библиотека `scikit-learn` располагает большим количеством моделей на основе метода ближайших соседей (<http://scikit-learn.org/stable/modules/neighbors.html>).

# Наивный Байес

Хорошо наивным быть для сердца, но вредно для ума.

*Анатоль Франс<sup>1</sup>*

От социальной сети мало толка, если люди в ней не могут общаться. Поэтому DataSciencester предлагает популярную среди пользователей соцсети возможность, которая позволяет им отправлять сообщения другим пользователям. И хотя большинство из них являются ответственными гражданами, которые отправляют только хорошо принимаемые сообщения, типа "как дела?", несколько злоумышленников настойчиво спамят других членов, рассылая незапрошенные адресатами сообщения по поводу схем быстрого обогащения, безрецептурной фармацевтической продукции и платных программ аккредитации в области науки о данных. Пользователи начали жаловаться, и поэтому директор по связям попросил вас при помощи методов науки о данных найти способ фильтрации спамных сообщений.

## Действительно глупый спам-фильтр

Пусть дано "вероятностное пространство", состоящее из множества всех возможных сообщений, и пусть приходит сообщение, выбранное из него случайным образом, где  $S$  — это событие "сообщение является спамом", а  $V$  — событие "сообщение содержит слово *виагра*". Тогда, согласно теореме Байеса, вероятность спамного сообщения со словом *виагра* равна:

$$P(S|V) = \frac{P(V|S)P(S)}{P(V|S)P(S) + P(V|\neg S)P(\neg S)},$$

где числитель обозначает условную вероятность, что сообщение является спамным и при этом содержит слово *виагра*, а знаменатель — безусловную вероятность встретить слово *виагра* во всем корпусе сообщений. Следовательно, эту формулу можно рассматривать, как способ представить долю спамных сообщений со словом *виагра*.

При наличии большого корпуса сообщений, о которых известно, что они спамные, и такого же корпуса сообщений, о которых известно, что они неспамные, можно легко оценить  $P(V|S)$  и  $P(V|\neg S)$ . Если далее допустить, что любое сообще-

<sup>1</sup> Анатоль Франс (1844–1924) — французский писатель и литературный критик. Лауреат Нобелевской премии по литературе (1921), деньги которой он пожертвовал в пользу голодающих России (см. [https://ru.wikipedia.org/wiki/Франс,\\_Анатоль](https://ru.wikipedia.org/wiki/Франс,_Анатоль)). — Прим. пер.

ние равновероятно является спамным или неспамным (благодаря чему  $P(S) = P(\neg S) = 0.5$ ), тогда:

$$P(S|V) = \frac{P(V|S)}{P(V|S) + P(V|\neg S)}.$$

Например, если слово *виагра* имеют 50% спамных сообщений и только 1% — неспамных, то вероятность, что любое конкретное сообщение со словом *виагра* является спамным, равна:

$$\frac{0.5}{0.5 + 0.01} = 90\%.$$

## Более продуманный спам-фильтр

Теперь пусть имеется словарь, состоящий из множества слов  $w_1, \dots, w_n$ . Чтобы перейти на язык теории вероятностей, обозначим  $X_i$  как событие "сообщение содержит слово  $w_i$ ". И пусть (благодаря некоему неопределенному на данный момент процессу) получена оценка  $P(X_i|S)$  вероятности, что спамное сообщение содержит  $i$ -е слово, а также аналогичная оценка  $P(X_i|\neg S)$  вероятности, что неспамное сообщение содержит  $i$ -е слово.

В основе *наивного байесовского классификатора* (naive Bayes classifier) лежит (серьезное) допущение, что наличие (или отсутствие) каждого конкретного слова не зависит от других слов в спамном или неспамном сообщении. Это допущение интуитивно подразумевает, что знание о том, что некое спамное сообщение содержит слово "виагра", не дает никакой информации о том, содержит ли то же самое сообщение слово "ролекс". Математически это свойство выражается следующим образом:

$$P(X_1 = x_1, \dots, X_n = x_n | S) = P(X_1 = x_1 | S) \cdots P(X_n = x_n | S).$$

Это крайне упрощающее допущение. (Оно объясняет, почему этот алгоритм называют *наивным*.) Представим, что словарь состоит *только* из слов "виагра" и "ролекс", и что половина всех спамных сообщений содержит выражение "дешевая виагра", а другая половина — "подлинные часы Rolex". В этом случае *наивная байесовская оценка вероятности, что спамное сообщение содержит слова "виагра" и "ролекс", будет следующей:*

$$P(X_1 = 1, X_2 = 1 | S) = P(X_1 = 1 | S) \cdot P(X_2 = 1 | S) = 0.5 \cdot 0.5 = 0.25,$$

поскольку эта оценка построена на предположении, что слова "виагра" и "ролекс" практически никогда не встречаются вместе. Несмотря на нереалистичность такого допущения, эта модель часто дает хорошие результаты и используется в реальных спам-фильтрах.

Кроме того, согласно тому же рассуждению на основе теоремы Байеса, которое использовалось для спам-фильтра, предназначенного только для *одного* слова *виагра*,



можно рассчитать вероятность спамности всего сообщения, используя для этого следующее равенство:

$$P(S | X = x) = \frac{P(X = x | S)}{P(X = x | S) + P(X = x | \neg S)}$$

Наивное байесовское допущение позволяет вычислить каждую из вероятностей справа, просто перемножив между собой индивидуальные вероятностные оценки для каждого слова словаря.

На практике обычно обходятся без перемножения большого количества вероятностей между собой во избежание проблемы арифметического *переполнения снизу* (приводящего к исчезновению разрядов), с которой компьютеры не справляются при работе с числами с плавающей точкой, находящимися слишком близко к нулю. Ссылаясь на свойства логарифмов из алгебры:  $\log ab = \log a + \log b$  и  $\exp(\log x) = x$ , произведение  $p_1 \cdots p_n$  обычно представляют в виде эквивалентной (и более дружественной к числам с плавающей точкой) суммы логарифмов:

$$\exp(\log(p_1) + \dots + \log(p_n)).$$

Единственная задача, которую осталось решить, — это получить оценки для  $P(X_i | S)$  и  $P(X_i | \neg S)$ , что спамное (или неспамное) сообщение содержит слово  $w_i$ . При наличии значительного числа "обучающих" сообщений, помеченных как спам и неспам, очевидное решение заключается в оценке  $P(X_i | S)$  просто в виде доли спамных сообщений со словом  $w_i$ .

Однако есть одно большое препятствие. Предположим, что в обучающей выборке слово "данные" присутствует только в неспамных сообщениях. Тогда оценка для  $P(\text{"данные"} | S)$  будет равна 0. В результате наивный байесовский классификатор всегда будет назначать вероятность спама, равную 0, *любому* сообщению со словом "данные", даже сообщениям типа "данные о дешевой виагре и подлинных часах Rolex". Во избежание этой проблемы обычно используют один из видов сглаживания.

В частности, воспользуемся константой сглаживания  $k$  и будем оценивать вероятность встретить  $i$ -е слово в спамном сообщении так:

$$P(X_i | S) = \frac{k + \text{число спамных сообщений с } w_i}{2k + \text{число спамных сообщений}}$$

Аналогичным образом для  $P(X_i | \neg S)$ . Другими словами, при вычислении вероятностей спама с  $i$ -м словом предполагается, что также встретилось  $k$  дополнительных спамных сообщений с этим словом и  $k$  дополнительных спамных сообщений без этого слова.

Например, если слово "данные" присутствует в 0/98 спамных документах обучающей выборки и  $k = 1$ , то оценка  $P(\text{"данные"} | S)$  будет  $1/100 = 0.01$ , что позволит

классификатору все равно назначать спамным сообщениям со словом "данные" ненулевые вероятности.

## Реализация

Теперь есть все составляющие, необходимые для построения классификатора. Прежде всего создадим простую функцию, которая, получая текст сообщения, генерирует список слов без повторов. Сначала преобразуем буквы каждого сообщения в строчные. Затем воспользуемся методом `re.findall()` для извлечения "слов", состоящих из букв, цифр и апострофов, и, наконец, функцией `set()` для получения списка уникальных слов:

```
# сгенерировать список слов без повторов
def tokenize(message):
    message = message.lower()           # преобразовать в строчные
    all_words = re.findall("[a-z0-9'+]", message) # извлечь слова
    return set(all_words)                # удалить повторы
```

Вторая функция подсчитывает частотность слов в маркированной обучающей выборке сообщений. Она возвращает словарь, ключами которого являются слова, а значения — двухэлементные списки `[spam_count, non_spam_count]`, соответствующие тому, сколько раз конкретное слово встречалось в спамных и неспамных сообщениях:

```
# подсчитать частотность слов
def count_words(training_set):
    """обучающая выборка состоит из пар (сообщение, спам?)"""
    counts = defaultdict(lambda: [0, 0])
    for message, is_spam in training_set:
        for word in tokenize(message):
            counts[word][0 if is_spam else 1] += 1
    return counts
```

На следующем шаге эти частотности преобразуются в оценки вероятностей, полученные при помощи метода сглаживания, описанного выше. Функция возвращает список триплетов, содержащих конкретное слово и значения вероятностей встретить его в спамном и неспамном сообщениях:

```
# вероятности слов
def word_probabilities(counts, total_spams, total_non_spams, k=0.5):
    """преобразовать частотности word_counts в список триплетов
    слово w, p(w | spam) и p(w | ~spam)"""
    return [(w,
              (spam + k) / (total_spams + 2 * k),
              (non_spam + k) / (total_non_spams + 2 * k))
            for w, (spam, non_spam) in counts.iteritems()]
```

В последнем фрагменте вероятности слов (и допущения наивного байесовского классификатора) используются для назначения вероятностей сообщениям:

```

# вероятность спама
def spam_probability(word_probs, message):
    message_words = tokenize(message)
    log_prob_if_spam = log_prob_if_not_spam = 0.0

    # просмотреть все слова в словаре
    for word, prob_if_spam, prob_if_not_spam in word_probs:

        # если СЛОВО в сообщении появляется, то
        # добавить лог-вероятность встретить его в сообщении
        if word in message_words:
            log_prob_if_spam += math.log(prob_if_spam)
            log_prob_if_not_spam += math.log(prob_if_not_spam)
        # если СЛОВО в сообщении не появляется, то
        # добавить лог-вероятность НЕ встретить его в сообщении,
        # вычисляемое как  $\log(1 - \text{вероятность встретить его в сообщении})$ 
        else:
            log_prob_if_spam += math.log(1.0 - prob_if_spam)
            log_prob_if_not_spam += math.log(1.0 - prob_if_not_spam)

    prob_if_spam = math.exp(log_prob_if_spam)
    prob_if_not_spam = math.exp(log_prob_if_not_spam)
    return prob_if_spam / (prob_if_spam + prob_if_not_spam)

```

Теперь можно собрать все вместе в питоновском классе NaiveBayesClassifier (*наивный байесовский классификатор*):

```

class NaiveBayesClassifier:

    def __init__(self, k=0.5):
        self.k = k
        self.word_probs = []

    # обучить классификатор при помощи обучающей выборки
    def train(self, training_set):

        # подсчитать спамные и неспамные сообщения
        num_spams = len([is_spam
                        for message, is_spam in training_set
                        if is_spam])
        num_non_spams = len(training_set) - num_spams

        # пропустить обучающую выборку через "конвейер"
        word_counts = count_words(training_set)
        self.word_probs = word_probabilities(word_counts,
                                           num_spams,
                                           num_non_spams,
                                           self.k)

```

```
def classify(self, message):
    return spam_probability(self.word_probs, message)
```

## Тестирование модели

Возьмем хороший (хотя и немного устаревший) набор данных под названием публичный корпус SpamAssassin (<https://spamassassin.apache.org/publiccorpus/>), в котором воспользуемся файлами с префиксом 20021010. (В Windows понадобится программа 7-zip (<http://www.7-zip.org/>), чтобы распаковать и извлечь их.)

После извлечения данных (скажем, в каталог C:\spam) должно получиться три папки: spam, easy\_ham и hard\_ham (спам, легкий\_неспам, трудный\_неспам). Каждая папка будет содержать набор писем, каждое из которых находится в отдельном файле. Чтобы не усложнять задачу, обратимся к строке с темой каждого письма.

Как определить, где находится тема? Просматривая файлы, можно обнаружить, что все они начинаются с последовательности "Subject:". Поэтому будем искать ее:

```
import glob, re

# поменяйте этот путь на путь фактического размещения файлов
path = r"C:\spam\*\*"

data = []

# glob.glob возвращает имена файлов,
# соответствующие шаблону поиска по указанному пути
for fn in glob.glob(path):
    is_spam = "ham" not in fn

    with open(fn, 'r') as file:
        for line in file:
            if line.startswith("Subject:"):
                # удалить начальное слово "Subject: " и взять все остальные
                subject = re.sub(r"^Subject: ", "", line).strip()
                data.append((subject, is_spam))
```

Теперь можно подразделить данные на обучающую и контрольную выборки, после чего все готово для построения модели классификатора:

```
random.seed(0) # чтобы получить повторяемые ответы
train_data, test_data = split_data(data, 0.75)
classifier = NaiveBayesClassifier() # инициализировать экземпляр класса
classifier.train(train_data) # обучить модель
```

и теперь можно проверить работу модели:

```
# триплеты (тема, is_spam по факту, предсказанная вероятность спама)
classified = [(subject, is_spam, classifier.classify(subject))
              for subject, is_spam in test_data]
```

```
# допустить, что spam_probability > 0.5 соответствует предсказанию спама,
# и подсчитать комбинации (спам фактический, спам предсказанный)
counts = Counter((is_spam, spam_probability > 0.5)
                 for _, is_spam, spam_probability in classified)
```

Это даст 101 истинноположительный (спам классифицируется как "спам"), 33 ложноположительных (неспам классифицируется как "спам"), 704 истинноотрицательных (неспам классифицируется как "не-спам") и 38 ложноотрицательных (спам классифицируется как "не-спам") результатов. Следовательно, точность и полнота равны соответственно  $101/(101+33) = 75\%$  и  $101/(101+38) = 73\%$ , что совсем неплохо для такой простой модели.

Кроме этого, интересно взглянуть на предельные случаи ошибочной классификации:

```
# отсортировать по вероятности спама от наименьшей до наибольшей
classified.sort(key=lambda row: row[2])
```

```
# максимальная предсказанная вероятность спама
# среди неспамных сообщений
spammiest_hams = filter(lambda row: not row[1], classified)[-5:]
```

```
# минимальная предсказанная вероятность спама
# среди фактически спамных сообщений
hammiest_spams = filter(lambda row: row[1], classified)[:5]
```

Два отнесенных к спаму неспамных сообщения содержат слова *needed* (с вероятностью появиться в спаме в 77 раз больше), *insurance* (с вероятностью появиться в спаме в 30 раз больше) и *important* (с вероятностью появиться в спаме в 10 раз больше).

Первое отнесенное к неспаму спамное сообщение слишком короткое ("Re: girls"), чтобы о чем-то судить; второе относится к вымогательству кредитной карты, большая часть слов которого отсутствует в обучающей выборке.

Точно так же можно взглянуть на самые спамные слова:

```
# вероятность спама при наличии заданного слова
def p_spam_given_word(word_prob):
    """использует теорему Байеса для вычисления вероятности
    p(спам | сообщение содержит слово)"""

    # word_prob - это один из триплетов,
    # созданный функцией word_probabilities
    # слово и его вероятности в случае спама и неспама
    word, prob_if_spam, prob_if_not_spam = word_prob
    return prob_if_spam / (prob_if_spam + prob_if_not_spam)
```

```
words = sorted(classifier.word_probs, key=p_spam_given_word)
```

```
spammiest_words = words[-5:]
hammiest_words = words[:5]
```

Самые спамные слова — money, systemworks, rates, sale и year, — по-видимому, связаны с попытками заставить людей покупать вещи. Большинство самых неспамных слов — spambayes, users, razor, zzzzteana и sadev, — как ни странно, похоже, связаны с предотвращением спама.

Каким образом можно улучшить производительность классификатора? Один из очевидных способов — предоставить больше обучающих данных. Кроме этого, есть ряд других способов улучшения модели. Вот несколько возможностей, которые стоит попробовать.

- ◆ Просматривать не только темы, но и содержание сообщений. Следует внимательно отнестись к обработке заголовков сообщений.
- ◆ Данный классификатор учитывает все слова, которые появляются в обучающей выборке, даже те, которые появляются всего один раз. Следует модифицировать классификатор так, чтобы он учитывал дополнительное пороговое значение `min_count` и игнорировал слова, число появлений которых меньше порогового значения.
- ◆ Генератор лексем не имеет понятия о похожих словах (к примеру, `shear` и `shearpest`). Можно модифицировать классификатор, включив дополнительную функцию морфологического анализа, которая будет преобразовывать слова в *классы эквивалентности*. Например, в качестве очень простого морфологического анализатора может выступать следующая функция:

```
# удалить окончание s
def drop_final_s(word):
    return re.sub("s$", "", word)
```

Создание хорошего морфологического анализатора — задача сложная. Для нахождения основы слова часто пользуются алгоритмом Портера (<http://tartarus.org/martin/PorterStemmer/>, [https://ru.wikipedia.org/wiki/Стеммер\\_Портера](https://ru.wikipedia.org/wiki/Стеммер_Портера)).

- ◆ Используемые признаки имеют одинаковую форму "сообщение содержит слово  $w_i$ ". Однако ничего не мешает добавить в текущую реализацию дополнительные признаки, такие как "сообщение содержит число", путем создания фиктивных слов, таких как *содержит:число*, и изменив генератор лексем `tokenizer` так, чтобы он порождал их в случае необходимости.

## Для дальнейшего изучения

- ◆ Статьи Пола Грэма "План для спама" (<http://www.paulgraham.com/spam.html>) и "Улучшенная байесовская фильтрация" (<http://www.paulgraham.com/better.html>) могут (быть интересными и) углубить понимание концепции, лежащей в основе построения спам-фильтров.
- ◆ Библиотека `scikit-learn` ([http://scikit-learn.org/stable/modules/naive\\_bayes.html](http://scikit-learn.org/stable/modules/naive_bayes.html)) содержит модель `BernoulliNB`, которая реализует тот же самый алгоритм наивной байесовской классификации, который реализован здесь, а также некоторые другие его версии, основанные на этой модели.

# Простая линейная регрессия

В морали, как в живописи, главное состоит в том, чтобы в нужном месте провести линию.

Г. К. Честертон<sup>1</sup>

В главе 5 для измерения силы линейной связи между двумя переменными использовалась функция корреляции *correlation*. Однако в большинстве задач недостаточно знать, что такая линейная связь существует. Необходимо иметь возможность установить природу связи. Именно для этих целей применяется модель простой (парной или однофакторной) линейной регрессии.

## Модель

Напомним, что в той главе исследовалась связь между числом друзей пользователя социальной сети DataSciencester и количеством времени, которое он проводит на ее сайте каждый день. Пусть мы убедили себя в том, что увеличение числа друзей *ведет к* тому, что приходится проводить на сайте больше времени, в отличие от альтернативных объяснений, которые обсуждались.

Директор по взаимодействию просит вас построить модель, описывающую эту связь. Поскольку вы нашли довольно сильную линейную зависимость, то логично сначала применить линейную модель.

В частности, можно предположить, что существуют такие коэффициенты  $\alpha$  (альфа) и  $\beta$  (бета), что:

$$y_i = \beta x_i + \alpha + \epsilon_i,$$

где  $y_i$  — это результирующий признак (или зависимая переменная), представляющий число минут, которые пользователь  $i$  ежедневно проводит на сайте;  $x_i$  — это факторный признак (или независимая переменная, или предиктор), представляющий число друзей пользователя  $i$ ;  $\epsilon_i$  — это случайная ошибка (в надежде небольшая) в силу других факторов, неучтенных в этой простой модели.

Пусть такие коэффициенты  $\alpha$  и  $\beta$  установлены. Тогда прогнозирование выполняется буквально следующим образом:

---

<sup>1</sup> Гилберт Кит Честертон (1874–1936) — английский христианский мыслитель, журналист и писатель (см. [https://ru.wikipedia.org/wiki/Честертон,\\_Гилберт\\_Кит](https://ru.wikipedia.org/wiki/Честертон,_Гилберт_Кит)). — *Прим. пер.*

```
# функция прогнозирования
def predict(alpha, beta, x_i):
    return beta * x_i + alpha
```

Как оценить  $\alpha$  и  $\beta$ ? Начнем с того, что результат прогнозирования для каждого входящего значения  $x_i$  факторного признака будет получен при любом выборе  $\alpha$  и  $\beta$ . И поскольку фактическое значение  $y_i$  результативного признака известно, то для каждой пары можно рассчитать ошибку, или отклонение, что то же самое:

```
# вернуть ошибку, т. е. отклонение наблюдаемого значения
# зависимой переменной y_i от модельных значений predict
def error(alpha, beta, x_i, y_i):
    """ошибка в результате прогнозирования beta * x_i + alpha
    при наблюдаемом значении y_i"""
    return y_i - predict(alpha, beta, x_i)
```

В первую очередь интересует суммарная ошибка по всему набору данных. Однако ошибки не просто суммируются — если прогноз для  $x_1$  слишком высокий, а для  $x_2$  — слишком низкий, то в результате ошибки могут взаимопогаситься.

Вместо этого суммируются *квадраты* ошибок:

```
# сумма квадратичных ошибок
def sum_of_squared_errors(alpha, beta, x, y):
    return sum(error(alpha, beta, x_i, y_i) ** 2
                for x_i, y_i in zip(x, y))
```

Для нахождения коэффициентов  $\alpha$  и  $\beta$ , таких, которые дают наименьшую сумму квадратов ошибок `sum_of_squared_errors`, используется *метод наименьших квадратов* (МНК, least squares).

Используя исчисление (или скучную алгебру), коэффициенты  $\alpha$  и  $\beta$ , минимизирующие величину ошибки, задаются при помощи:

```
# подгонка методом наименьших квадратов
def least_squares_fit(x, y):
    """при заданных обучающих значениях x и y
    найти значения alpha и beta на основе МНК"""
    beta = correlation(x, y) * standard_deviation(y) / standard_deviation(x)
    alpha = mean(y) - beta * mean(x)
    return alpha, beta
```

Не углубляясь в точную математику, подумаем, почему это может быть разумным решением. Выбор коэффициента  $\alpha$  всего лишь говорит, что при наличии среднего значения независимой переменной  $x$  мы прогнозируем среднее значение зависимой переменной  $y$ .

Выбор коэффициента  $\beta$  означает, что при увеличении входящего значения на стандартное отклонение `standard_deviation(x)` прогнозное значение увеличивается на `correlation(x, y) * standard_deviation(y)`. В случае, когда  $x$  и  $y$  идеально коррелированы, увеличение на единицу стандартного отклонения в  $x$  приводит к росту на



единицу стандартного отклонения  $y$  в прогнозе. Когда они идеально антикоррелированы, увеличение в  $x$  ведет к *снижению* в прогнозе. И когда корреляция нулевая, то  $\beta = 0$ , и, следовательно, изменения в  $x$  совершенно не влияют на прогноз.

Эту функцию несложно применить к данным без выбросов из главы 5:

```
alpha, beta = least_squares_fit(num_friends_good, daily_minutes_good)
```

Функция вернет значения  $\alpha = 22.95$  и  $\beta = 0.903$ . Следовательно, согласно модели, ожидается, что пользователь с  $n$  друзьями из соцсети DataSciencester будет проводить на сайте  $22.95 + n * 0.903$  минут ежедневно. Другими словами, модель предсказывает, что пользователь без друзей все равно будет проводить на сайте около 23 минут в день, а каждый дополнительный друг, как ожидается, будет увеличивать проводимое пользователем время почти на одну минуту в день.

На рис. 14.1 показана прогнозная линия, которая дает представление о качестве подгонки модели к наблюдаемым данным.

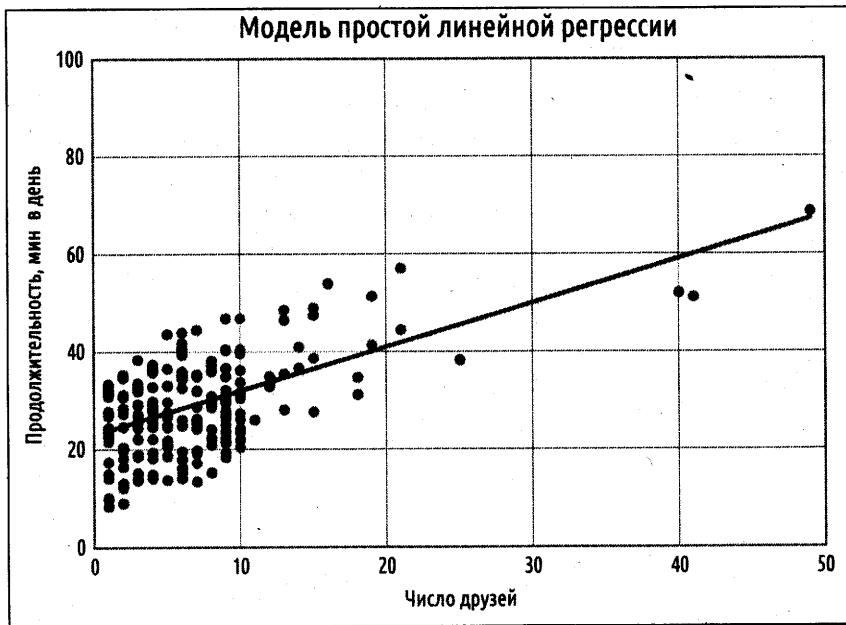


Рис. 14.1. Модель простой линейной регрессии

Разумеется, для того чтобы выяснить качество подгонки к данным, требуется более подходящий способ, чем просто рассматривать диаграмму. Распространенной мерой оценки качества подгонки является *коэффициент детерминации* (или *R-квадрат*), который измеряет долю суммарной вариации в зависимой переменной, объясняемой моделью:

```
# полная сумма квадратов
def total_sum_of_squares(y):
    """полная сумма квадратов отклонений  $y_i$  от их среднего"""
    return sum(v ** 2 for v in de_mean(y))
```

# R-квадрат — это доля вариации зависимой переменной, объясняемая моделью

```
def r_squared(alpha, beta, x, y):
```

```
    """доля вариации в y, объясненная моделью, равна
    1 - доля вариации в y, не объясненная моделью"""
```

```
    return 1.0 - (sum_of_squared_errors(alpha, beta, x, y) /
                  total_sum_of_squares(y))
```

```
r_squared(alpha, beta, num_friends_good, daily_minutes_good) # 0.329
```

Итак, выбраны коэффициенты  $\alpha$  и  $\beta$ , которые минимизируют сумму квадратичных ошибок предсказания. Однако можно было бы выбрать линейную модель, которая "всегда предсказывает среднее значение  $\text{mean}(y)$ " (что соответствует  $\alpha = \text{mean}(y)$  и  $\beta = 0$ ), чья сумма квадратичных ошибок в точности равна полной сумме квадратов. И тогда R-квадрат был бы равен нулю, указывая на модель, которая (со всей очевидностью) работает не лучше, чем простой прогноз среднего.

Ясно, что модель на основе МНК должна как минимум прогнозировать не хуже, и поэтому сумма квадратичных ошибок должна быть *не больше* полной суммы квадратов, а значит, R-квадрат должен быть не меньше нуля. С другой стороны, сумма квадратичных ошибок должна быть не меньше 0, вследствие чего R-квадрат может быть не больше 1.

Чем ближе число к 1, тем лучше подгонка модели к данным. Здесь R-квадрат равен 0.329, что говорит о том, что модель соответствует данным только частично, и что очевидно присутствуют дополнительные факторы.

## Применение метода градиентного спуска

Если записать  $\theta = [\alpha, \beta]$ , то эту задачу можно решить методом градиентного спуска:

```
# квадратичная ошибка
```

```
def squared_error(x_i, y_i, theta):
```

```
    alpha, beta = theta
```

```
    return error(alpha, beta, x_i, y_i) ** 2
```

```
# градиент квадратичной ошибки
```

```
def squared_error_gradient(x_i, y_i, theta):
```

```
    alpha, beta = theta
```

```
    # частные производные alpha и beta
```

```
    return [-2 * error(alpha, beta, x_i, y_i),
            -2 * error(alpha, beta, x_i, y_i) * x_i]
```

```
# выбрать случайное значение для запуска
```

```
random.seed(0)
```

```
theta = [random.random(), random.random()]
```

```
alpha, beta = minimize_stochastic(squared_error,
                                  squared_error_gradient,
                                  num_friends_good,
                                  daily_minutes_good,
                                  theta,
                                  0.0001)

print(alpha, beta)
```

Пользуясь теми же данными, получаем  $\alpha = 22.93$  и  $\beta = 0.905$ , которые очень близки к точным ответам.

## Метод максимального правдоподобия

Почему используется метод наименьших квадратов? Одно из объяснений касается оценивания неизвестного параметра по методу *максимального правдоподобия* (maximum likelihood estimation)<sup>2</sup>.

Пусть имеется выборка данных  $v_1, \dots, v_n$ , которая порождается распределением, зависящим от некоторого неизвестного параметра  $\theta$ :

$$p(v_1, \dots, v_n | \theta).$$

Если параметр  $\theta$  (тэта) неизвестен, то можно поменять члены местами, чтобы представить эту величину, как *правдоподобие* параметра  $\theta$  при наличии выборки:

$$L(\theta | v_1, \dots, v_n).$$

Согласно такому подходу наиболее правдоподобным значением  $\theta$  является то, которое максимизирует эту функцию правдоподобия, т. е. значение, которое делает наблюдаемые данные наиболее вероятными. В случае непрерывного распределения, где вместо функции массы вероятности используется функция распределения вероятностей<sup>3</sup>, можно сделать то же самое.

Но вернемся к регрессионному анализу. Одно из допущений, которое нередко принимается относительно модели простой линейной регрессии, заключается в том, что случайные ошибки регрессии нормально распределены с нулевым средним и неким (известным) стандартным отклонением  $\sigma$ . В этом случае правдоподобие на основе наблюдаемой пары  $(x_i, y_i)$  равно:

$$L(\alpha, \beta | x_i, y_i, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-(y_i - \alpha - \beta x_i)^2}{2\sigma^2}\right).$$

<sup>2</sup> Оценка максимального правдоподобия — это метод оценивания неизвестного параметра путем максимизации функции вероятности, в результате которой приобретаются значения параметров модели, которые делают данные "ближе" к реальным. — *Прим. пер.*

<sup>3</sup> В зависимости от контекста под *функцией распределения вероятностей* подразумевается либо дифференциальная функция распределения (ДФР, плотность распределения), либо интегральная функция распределения (ИФР, кумулятивная функция). Функция массы вероятности (probability mass function, PMF) дает относительную вероятность того, что случайная дискретная величина в точности равна некоторому значению. — *Прим. пер.*

Правдоподобие на основе всего набора данных является произведением индивидуальных правдоподобий, которое максимально именно тогда, когда полученные оценочные значения  $\alpha$  и  $\beta$  дают минимальную сумму квадратов случайных ошибок. Другими словами, в данном случае (и с этими допущениями) минимизация суммы квадратов ошибок эквивалентна максимизации правдоподобия наблюдаемых данных.

## **Для дальнейшего изучения**

Продолжайте знакомиться с методами регрессионного анализа! В *главе 15* речь пойдет о множественной регрессии.

# Множественная регрессия

Я не берусь за решение проблемы, не вкладывая в нее переменные, которые не смогут на нее повлиять.

*Билл Парцелл<sup>1</sup>*

Хотя директор порядком впечатлен вашей прогнозной моделью, он все же считает, что ее можно улучшить. С этой целью вы собрали дополнительные данные: по каждому пользователю теперь известно, сколько часов он работает каждый день и есть ли у него ученая степень. Вы собираетесь их использовать для усовершенствования модели.

В соответствии с этим вы строите гипотезу о линейной модели, но теперь уже с несколькими независимыми переменными или факторными признаками:

$$\text{минуты} = \alpha + \beta_1 \text{друзья} + \beta_2 \text{продолжительность} + \beta_3 \text{степень} + \varepsilon.$$

Очевидно, наличие у пользователя ученой степени не является числом, но, как уже упоминалось в *главе 11*, для придания количественной определенности можно ввести *фиктивную переменную*, которая будет равна 1 для пользователей с ученой степенью и 0 — для пользователей без нее, после чего она становится такой же численной переменной, как и другие.

## Модель

Напомним, что в *главе 14* выполнялся подбор модели следующего вида:

$$y_i = \alpha + \beta x_i + \varepsilon_i.$$

Теперь представим, что каждый входящий факторный признак  $x_i$  — это не одно число, а вектор из  $k$  чисел  $x_{i1}, \dots, x_{ik}$ . Тогда уравнение множественной (или многофакторной) модели регрессии будет следующим:

$$y_i = \alpha + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i,$$

где вектор параметров обозначается как  $\beta$ . В уравнении также имеется константный член<sup>2</sup>, и чтобы его учесть, добавим к исходной матрице данных единичный столбец:

<sup>1</sup> Дуэйн Чарльз "Билл" Парцеллс (1941) — бывший главный тренер американской национальной команды по американскому футболу (см. [https://ru.wikipedia.org/wiki/Веннингтон,\\_Билл](https://ru.wikipedia.org/wiki/Веннингтон,_Билл)). — *Прим. пер.*

```
beta = [alpha, beta_1, ..., beta_k]
```

и:

```
x_i = [1, x_i1, ..., x_ik]
```

Тогда прогнозная модель будет следующей:

```
def predict(x_i, beta):
    """предполагается, что первый элемент каждого x_i равен 1"""
    return dot(x_i, beta)
```

В данном конкретном случае независимая переменная  $x$  будет списком векторов, каждый из которых выглядит следующим образом:

```
[1,      # константа
 49,     # число друзей
 4,      # продолжительность рабочего времени в день, часы
 0]      # не имеет степени кандидата наук
```

## Другие допущения модели наименьших квадратов

Для того чтобы рассматриваемая модель и решение имели смысл, следует принять еще ряд допущений.

Первое заключается в том, что столбцы  $x$  *линейно независимы*, т. е. невозможно записать любой из них в виде взвешенной суммы каких-то других столбцов. Если это допущение не удовлетворяется, то оценить  $\beta$  невозможно. Чтобы убедиться, воспользуемся предельным случаем и представим, что в данных имелось дополнительное поле для количества знакомых `num_acquaintances`, которое для каждого пользователя в точности равно числу друзей `num_friends`.

Затем, начиная с любого  $\beta$ , если добавить к коэффициенту `num_friends` *любое* число и вычесть это же число из коэффициента `num_acquaintances`, то предсказания модели не изменятся. И, следовательно, найти значение коэффициента для `num_friends` не представляется возможным. (Обычно нарушения этого допущения не столь очевидны.)

Второе важное допущение касается того, что все столбцы  $x$  некоррелированы со случайными ошибками  $\varepsilon$ . Если это допущение не выполняется, то оценки  $\beta$  будут систематически неверными.

Например, в *главе 14* была построена модель, которая предсказывала, что каждый дополнительный друг связан с лишними 0.90 минутами, проводимыми на сайте ежедневно.

<sup>2</sup> Имеется в виду *пересечение* (англ. *intercept*) — свободный коэффициент, которому равна зависимая переменная, если предиктор, или независимая переменная, равен нулю. Геометрически, это точка, в которой линия регрессии пересекает ось  $y$ . Отсюда термин для этого коэффициента — *пересечение*. — *Прим. пер.*

Представим также, что.

- ◆ лица, которые работают больше часов, проводят на сайте меньше времени;
- ◆ лица, у которых больше друзей, как правило, работают больше часов.

То есть "фактическая" модель следующая:

$$\text{минуты} = \alpha + \beta_1 \text{друзья} + \beta_2 \text{продолжительность} + \varepsilon,$$

и что число друзей и продолжительность рабочего времени положительно коррелируют. В этом случае при минимизации ошибок однофакторной модели:

$$\text{минуты} = \alpha + \beta_1 \text{друзья} + \varepsilon$$

будет недооценен  $\beta_1$ .

Теперь подумаем, что произойдет, если бы предсказания выполнялись на основе однофакторной модели с таким "фактическим" значением  $\beta_1$ . (То есть значением, получаемом в результате минимизации ошибок модели, которую мы назвали "фактической".) Прогнозы будут, как правило, слишком малыми для пользователей, которые работают много часов, и слишком большими для пользователей, которые работают всего несколько часов, потому что  $\beta_2 > 0$ , а этот фактор "забыли" учесть. Поскольку продолжительность рабочего времени положительно коррелирует с числом друзей, то, значит, прогнозные значения, как правило, будут слишком малыми для пользователей с большим числом друзей и слишком большими для пользователей лишь с несколькими друзьями.

В результате можно уменьшить ошибки (однофакторной модели), уменьшив оценку параметра  $\beta_1$ , а значит  $\beta_1$ , который минимизирует ошибки, окажется меньше "фактического" значения. Другими словами, в данном случае однофакторное решение на основе МНК смещено в сторону недооценки  $\beta_1$ . И вообще, всегда, когда независимые переменные коррелируют с ошибками описанным образом, решение на основе МНК дает смещенную оценку  $\beta$ .

## Подбор модели

Так же как и в простой линейной модели, выполним оценку  $\beta$ , дающую минимальное значение суммы квадратов ошибок. Найти точное решение этой задачи вручную нелегко, поэтому воспользуемся методом градиентного спуска. Начнем с создания функции ошибок, которую требуется минимизировать. Для расчетов методом стохастического градиентного спуска нужно лишь, чтобы квадратичная ошибка соответствовала одиночному прогнозу:

```
# функция ошибок
def error(x_i, y_i, beta):
    return y_i - predict(x_i, beta)

# квадратичная ошибка
def squared_error(x_i, y_i, beta):
    return error(x_i, y_i, beta) ** 2
```

Читатели, разбирающиеся в исчислениях, могут вычислить:

```
# градиент квадратичной ошибки
def squared_error_gradient(x_i, y_i, beta):
    """градиент (относительно beta), соответствующий i-му квадрату ошибки"""
    return [-2 * x_ij * error(x_i, y_i, beta)
            for x_ij in x_i]
```

В противном случае, придется поверить на слово.

На данный момент все готово, чтобы найти оптимальное значение beta с использованием метода стохастического градиентного спуска:

```
# оценить beta
def estimate_beta(x, y):
    beta_initial = [random.random() for x_i in x[0]]
    return minimize_stochastic(squared_error,
                               squared_error_gradient,
                               x, y,
                               beta_initial,
                               0.001)

random.seed(0)
beta = estimate_beta(x, daily_minutes_good) # [30.63, 0.972, -1.868, 0.911]
```

В результате получится следующая модель:

минуты = 30.63 + 0.972 друзья – 1.868 продолжительность + 0.911 степень.

## Интерпретация модели

Коэффициенты модели следует рассматривать как оценки влияния каждого конкретного фактора в условиях, когда остальные факторы постоянны. При прочих равных каждый дополнительный друг соответствует лишней минуте, проводимой на сайте ежедневно. При прочих равных каждый дополнительный час в продолжительности рабочего времени соответствует сокращению времени, проводимого на сайте ежедневно, примерно на 2 минуты. При прочих равных наличие ученой степени связано с проведением на сайте ежедневно одной лишней минуты.

Однако модель (непосредственно) ничего не говорит о взаимодействии переменных между собой. Вполне возможно, что влияние фактора продолжительности рабочего времени неодинаково для людей с большим числом друзей и с малым числом друзей. Данная модель не объясняет это. Учесть этот случай можно, если ввести новую переменную, которая представляет собой *произведение* коэффициентов "число друзей" и "продолжительность рабочего времени", что фактически позволяет увеличивать (или уменьшать) коэффициент "продолжительность рабочего времени" по мере увеличения числа друзей.

Еще одна возможность заключается в том, что зависимость количества проводимого на сайте времени от числа друзей работает *до определенного момента*, после которого дальнейшее увеличение числа друзей ведет к уменьшению проводимого



на сайте времени. (Может быть, при слишком большом количестве друзей опыт взаимодействия станет обременительным?) Можно было бы попробовать объяснить это в модели, добавив еще одну переменную в виде *квадрата* числа друзей.

При этом, начиная добавлять переменные, следует позаботиться о том, чтобы они имели "смысл". Добавлять же произведения, логарифмы, квадраты и более высокие степени можно без ограничений.

## Качество подбора модели

Снова обратимся к коэффициенту детерминации (R-квадрату), который теперь увеличился до 0.68:

```
# многофакторный R-квадрат
def multiple_r_squared(x, y, beta):
    sum_of_squared_errors = sum(error(x_i, y_i, beta) ** 2
                                for x_i, y_i in zip(x, y))
    return 1.0 - sum_of_squared_errors / total_sum_of_squares(y)
```

Следует, однако, иметь в виду, что добавление к регрессии новых переменных *неизбежно* увеличивает R-квадрат. В конце концов, однофакторная модель регрессии является лишь частным случаем многофакторной модели, где коэффициенты "продолжительности рабочего времени" и "ученой степени" равны 0. Оптимальная модель многофакторной регрессии будет неизбежно иметь случайные ошибки как минимум столь же малые, что и у однофакторной.

По этой причине в модели множественной регрессии следует также обратиться к *стандартным ошибкам* коэффициентов, которые измеряют уровень уверенности модели в оценке каждого  $\beta$ . Линия регрессии в целом может очень хорошо соответствовать наблюдаемым данным, однако если некоторые независимые переменные коррелируют (или нерелевантны), то их коэффициенты могут не иметь какого-то *смысла*.

Типичный подход к измерению стандартных ошибок начинается с принятия еще одного допущения — о том, что случайные ошибки  $\varepsilon$ , представляют собой независимые нормальные случайные величины с нулевым средним значением и неким совместным (неизвестным) стандартным отклонением  $\sigma$ . В таком случае для нахождения стандартной ошибки каждого коэффициента можно воспользоваться линейной алгеброй (хотя скорее всего это сделает статистическое программное обеспечение). Чем ошибка больше, тем меньше наша модель уверена в данном коэффициенте. Увы, мы не готовы выполнить подобного рода линейные алгебраические вычисления с нуля.

## Отступление: бутстрапирование данных

Пусть имеется выборка из  $n$  точек данных, порожденная некоторым (неизвестным) распределением:

```
data = get_sample(num_points=n) # получить выборку из n точек
```

В главе 5 была реализована функция для вычисления медианы `median` наблюдаемых данных. Эту функцию можно использовать для оценки медианы непосредственно распределения.

Но насколько можно доверять этой оценке? Если все значения в выборке очень близки к 100, то скорее всего фактическая медиана тоже будет близка к 100. Если же примерно половина значений в выборке близка к 0, а другая половина близка к 200, то нельзя быть настолько же уверенным в отношении медианы.

Если бы имелась возможность неоднократно получать новые выборки, тогда можно было бы вычислить медиану каждой выборки и обратиться к распределению этих медиан. Обычно такая возможность отсутствует. Вместо этого пользуются *бутстрапированием* (*bootstrapping*) или методом размножения выборок, когда из имеющихся данных генерируют новые случайные выборки *с возвратом* размером *n* точек, а затем вычисляют медианы этих синтетических наборов данных<sup>3</sup>:

```
# выборка на основе бутстрап-метода
def bootstrap_sample(data):
    """случайно отбирает len(data) элементов с возвратом (с повторами)"""
    return [random.choice(data) for _ in data]

# статистика на основе бутстрап-метода
def bootstrap_statistic(data, stats_fn, num_samples):
    """оценивает функцию stats_fn на бутстрап-выборках из данных
    в количестве num_samples"""
    return [stats_fn(bootstrap_sample(data))
            for _ in range(num_samples)]
```

Например, рассмотрим следующие два набора данных:

```
# значения всех 101 точки очень близки к 100
close_to_100 = [99.5 + random.random() for _ in range(101)]

# значения 50 из 101 близки к 0, значения других 50 находятся рядом с 200
far_from_100 = ([99.5 + random.random()] +
                [random.random() for _ in range(50)] +
                [200 + random.random() for _ in range(50)])
```

Если вычислить медиану `median` каждой выборки, то обе будут находиться очень близко к 100. Однако если обратиться к:

```
bootstrap_statistic(close_to_100, median, 100)
```

то в большинстве случаев значение чисел будет действительно близко к 100, а если обратиться к:

```
bootstrap_statistic(far_from_100, median, 100)
```

то получим большое количество чисел, близких как к 0, так и 200.

<sup>3</sup> Выборка с возвратом (или выборка, полученная повторным способом) аналогична неоднократному вытягиванию случайной карты из колоды игральных карт, когда после каждого вытягивания карта возвращается назад в колоду. В результате время от времени обязательно будет попадаться карта, которая уже выбиралась. — *Прим. пер.*

Стандартное отклонение `standard_deviation` первого набора медиан имеет значение, близкое к 0, тогда как этот показатель для второго набора медиан имеет значение, близкое к 100. (В этом предельном случае достаточно легко убедиться, проверив данные вручну, хотя обычно так не поступают.)

## Стандартные ошибки коэффициентов регрессии

Аналогичный подход может быть принят относительно оценки стандартных ошибок коэффициентов регрессии. При помощи бутстрап-функции `bootstrap_sample` из данных неоднократно делается случайная выборка и на ее основе оценивается `beta`. Если коэффициент, соответствующий одной из независимых переменных (скажем, `num_friends`), варьирует по выборкам не сильно, то можно быть уверенными, что оценка близка к оптимальной. Если же коэффициент от выборки к выборке сильно варьирует, то уверенность в оценке минимальна.

Единственная тонкость состоит в том, что перед выполнением выборки необходимо, пользуясь функцией `zip`, объединить данные двух списков `x` и `y` в список кортежей с тем, чтобы обеспечить совместный отбор соответствующих значений независимой и зависимой переменных. Поэтому функция `bootstrap_sample` будет возвращать список пар `(x_i, y_i)`, которые потом нужно снова собрать в списки `x_sample` и `y_sample`:

```
# оценить beta для выборки
def estimate_sample_beta(sample):
    """выборка представлена списком пар (x_i, y_i)"""
    x_sample, y_sample = zip(*sample) # трюк с "разъединением" списка
    return estimate_beta(x_sample, y_sample)

random.seed(0) # благодаря этому можно получить повторяемые результаты

bootstrap_betas = bootstrap_statistic(zip(x, daily_minutes_good),
                                     estimate_sample_beta,
                                     100)
```

Затем можно оценить стандартное отклонение каждого коэффициента:

```
# стандартные ошибки сгенерированных выборок
bootstrap_standard_errors = [
    standard_deviation([beta[i] for beta in bootstrap_betas])
    for i in range(4)]

# [1.174, # константа, фактическая ошибка = 1.19
# 0.079, # число друзей, фактическая ошибка = 0.080
# 0.131, # не занят, фактическая ошибка = 0.127
# 0.990] # степень, фактическая ошибка = 0.998
```

Эти данные можно использовать для проверки статистических гипотез, таких как гипотеза о равенстве  $\beta_1$  нулю. При нулевой гипотезе  $\beta_1 = 0$  (и с другими допущениями о распределении  $\varepsilon_i$ ) статистика:

$$t_j = \hat{\beta}_j / \hat{\sigma}_j,$$

которая представляет собой оценку коэффициента  $\beta_1$ , деленную на оценку его стандартной ошибки, подчиняется *t-распределению Стьюдента*<sup>4</sup> с " $n - k$  степенями свободы".

Если бы имелась ИФР для *t-распределения Стьюдента* `students_t_cdf`, то можно было бы вычислить *p-значения* для каждого коэффициента по МНК, чтобы показать уровень правдоподобия наблюдать такое значение в случае, если фактический коэффициент равен нулю. К сожалению, такая функция отсутствует. (Хотя ее можно было бы реализовать, если бы работа велась не с чистого листа.)

Тем не менее, по мере увеличения степеней свободы *t-распределение* стремится к стандартному нормальному распределению. В подобной ситуации, когда  $n$  намного больше  $k$ , можно воспользоваться интегральной функцией для нормального распределения `normal_cdf` и при этом остаться довольным собой:

```
# p-значение
def p_value(beta_hat_j, sigma_hat_j):
    if beta_hat_j > 0:
        # если коэффициент имеет положительное значение,
        # нужно дважды вычислить вероятность наблюдать
        # еще более КРУПНОЕ значение
        return 2 * (1 - normal_cdf(beta_hat_j / sigma_hat_j))
    else:
        # в противном случае вероятность наблюдать МЕНЬШЕЕ значение
        return 2 * normal_cdf(beta_hat_j / sigma_hat_j)
```

```
p_value(30.63, 1.174) # ~0 (постоянный терм)
p_value(0.972, 0.079) # ~0 (количество друзей)
p_value(-1.868, 0.131) # ~0 (часы работы)
p_value(0.911, 0.990) # 0.36 (степень)
```

(В других условиях, вероятно, было бы использовано статистическое программное обеспечение, которое способно вычислить и *t-распределение*, и точные стандартные ошибки.)

Пока большинство коэффициентов имеют очень малые *p-значения* (в предположении, что они действительно ненулевые), коэффициент для "ученой степени" не

<sup>4</sup> Распределение Стьюдента — это однопараметрическое семейство абсолютно непрерывных распределений, которое по сути представляет собой сумму нескольких нормально распределенных случайных величин. Чем больше величин, тем больше вероятность, что их сумма будет иметь нормальное распределение. Таким образом, количество суммируемых величин определяет важнейший параметр формы данного распределения — число степеней свободы (см. [https://ru.wikipedia.org/wiki/Распределение\\_Стьюдента](https://ru.wikipedia.org/wiki/Распределение_Стьюдента)). — Прим. пер.

"значимо" отличается от нуля, тем самым скорее повышая вероятность его случайности, чем информативности.

В более сложных сценариях регрессии иногда требуется проверить более сложные гипотезы о данных, как например, о неравенстве нулю как минимум одного из  $\beta_j$  или о равенстве  $\beta_1 = \beta_2$  и  $\beta_3 = \beta_4$ , которые можно выполнить при помощи *F-теста*<sup>5</sup>, но этот вопрос, увы, выходит за рамки данной книги.

## Регуляризация

На практике часто возникает потребность в применении линейной регрессии к наборам данных с большим числом переменных. Это создает пару лишних затруднений. Во-первых, чем больше переменных, тем больше шансов переобучить модель обучающей выборкой. И во-вторых, чем больше ненулевых коэффициентов, тем труднее в них разобраться. Если задача заключается в *объяснении* некоторого явления, то разреженная модель с тремя факторами в практическом плане может оказаться полезнее, чем модель чуть получше, но с сотнями факторов.

*Регуляризация* — это метод, при котором (с целью предотвращения переобучения) к случайным ошибкам регрессии добавляется штраф, который увеличивается с увеличением  $\beta$ . Затем это сочетание ошибок и штрафа минимизируется. Чем большее значение придается штрафу, тем больше наказываются крупные коэффициенты.

Например, в *гребневой регрессии* добавляется штраф, пропорциональный сумме квадратов  $\beta_{i-1}$ . ( $\beta_0$ , как правило, не штрафуются.)

```
# штраф в гребневой регрессии
# alpha - это ГИПЕРПАРАМЕТР, контролирующий, насколько жестким будет штраф
# иногда его называют лямбда-параметром, но этот термин уже имеет свое
# значение в языке Python
def ridge_penalty(beta, alpha):
    return alpha * dot(beta[1:], beta[1:])

# квадратичная ошибка гребневой регрессии
def squared_error_ridge(x_i, y_i, beta, alpha):
    """оценить ошибку плюс штраф для beta"""
    return error(x_i, y_i, beta) ** 2 + ridge_penalty(beta, alpha)
```

Гребневую регрессию затем можно подключить к методу градиентного спуска обычным образом:

```
def ridge_penalty_gradient(beta, alpha):
    """градиент только гребневого штрафа"""
```

<sup>5</sup> *F-тестом* или критерием Фишера называют любую проверку статистической гипотезы, статистика которой при выполнении нулевой гипотезы имеет распределение Фишера (*F-распределение*). — *Прим. пер.*

```

return [0] + [2 * alpha * beta_j for beta_j in beta[1:]]

def squared_error_ridge_gradient(x_i, y_i, beta, alpha):
    """градиент, соответствующий i-й квадратичной ошибке,
    включая гребневой штраф"""
    return vector_add(squared_error_gradient(x_i, y_i, beta),
                      ridge_penalty_gradient(beta, alpha))

# оценить beta гребневой регрессии
def estimate_beta_ridge(x, y, alpha):
    """применить градиентный спуск для подгонки гребневой регрессии
    со штрафом alpha"""
    beta_initial = [random.random() for x_i in x[0]]
    return minimize_stochastic(partial(squared_error_ridge, alpha=alpha),
                               partial(squared_error_ridge_gradient,
                                       alpha=alpha),
                               x, y,
                               beta_initial,
                               0.001)

```

При  $\alpha = 0$  штраф отсутствует совсем, и получаются те же самые результаты, что и раньше:

```

random.seed(0)
beta_0 = estimate_beta_ridge(x, daily_minutes_good, alpha=0.0)
# [30.6, 0.97, -1.87, 0.91]
dot(beta_0[1:], beta_0[1:]) # 5.26
multiple_r_squared(x, daily_minutes_good, beta_0) # 0.680

```

По мере увеличения  $\alpha$  качество подбора (функции регрессии) ухудшается, но размер  $\beta$  становится меньше:

```

beta_0_01 = estimate_beta_ridge(x, daily_minutes_good, alpha=0.01)
# [30.6, 0.97, -1.86, 0.89]
dot(beta_0_01[1:], beta_0_01[1:]) # 5.19
multiple_r_squared(x, daily_minutes_good, beta_0_01) # 0.680

```

```

beta_0_1 = estimate_beta_ridge(x, daily_minutes_good, alpha=0.1)
# [30.8, 0.95, -1.84, 0.54]
dot(beta_0_1[1:], beta_0_1[1:]) # 4.60
multiple_r_squared(x, daily_minutes_good, beta_0_1) # 0.680

```

```

beta_1 = estimate_beta_ridge(x, daily_minutes_good, alpha=1)
# [30.7, 0.90, -1.69, 0.085]
dot(beta_1[1:], beta_1[1:]) # 3.69
multiple_r_squared(x, daily_minutes_good, beta_1) # 0.676

```

```

beta_10 = estimate_beta_ridge(x, daily_minutes_good, alpha=10)
# [28.3, 0.72, -0.91, -0.017]

```

```
dot(beta_10[1:], beta_10[1:])      # 1.36
multiple_r_squared(x, daily_minutes_good, beta_10)    # 0.573
```

В данном случае коэффициент "ученая степень" исчезает при увеличении штрафа, что согласуется с предыдущим результатом, который значимо не отличался от нуля.



Обычно перед тем, как применить этот подход, вы сначала выполняете многомерное шкалирование данных при помощи функции `rescale`. Ведь, если поменять годовой стаж на вековой, то его коэффициент по МНК увеличится в 100 раз и сразу же будет оштрафован гораздо сильнее, несмотря на то, что модель одна и та же.

Еще один подход — это *лассо*-регрессия, в которой используется штраф:

```
# штраф лассо-регрессии
def lasso_penalty(beta, alpha):
    return alpha * sum(abs(beta_i) for beta_i in beta[1:])
```

В отличие от штрафа, применяемого в гребневой регрессии, который в целом коэффициенты уменьшает, лассо-штраф стремится обнулить их, что делает этот метод пригодным для обучения разреженных моделей. К сожалению, он не поддается обработке методом градиентного спуска, и поэтому решить с чистого листа задачу с его помощью не получится.

## Для дальнейшего изучения

- ◆ Регрессионный анализ базируется на богатой и обширной теории. Это еще одна тема, которая нуждается в изучении на основе учебника или, по крайней мере, большого числа статей из Википедии.
- ◆ Библиотека `scikit-learn` включает модуль `linear_model` ([http://scikit-learn.org/stable/modules/linear\\_model.html](http://scikit-learn.org/stable/modules/linear_model.html)), который предлагает модель линейной регрессии `LinearRegression`, аналогичную рассматриваемой в книге, а также методы гребневой регрессии (`Ridge`), лассо-регрессии (`Lasso`) и другие виды регуляризации.
- ◆ `Statsmodels` (<http://statsmodels.sourceforge.net/>) — это еще один модуль Python, который (среди всего прочего) содержит модели линейной регрессии.

# Логистическая регрессия

Многие говорят, что между гениальностью и сумасшествием прочерчена тонкая линия.

Не думаю, что она тонкая. На самом деле она представляет собой зияющую пропасть.

*Билл Бэйли<sup>1</sup>*

В главе 1 была кратко рассмотрена задача с попыткой предсказать, какие из пользователей DataSciencester оплачивают свои привилегированные аккаунты. В данной главе эта задача будет рассмотрена повторно.

## Задача

Пусть имеется набор анонимных данных примерно на 200 пользователей, содержащий зарплату каждого пользователя, его опыт работы в качестве аналитика данных и состояние его платежей за привилегированный аккаунт (рис. 16.1). Как обычно, когда в качестве зависимой переменной используются категориальные дан-

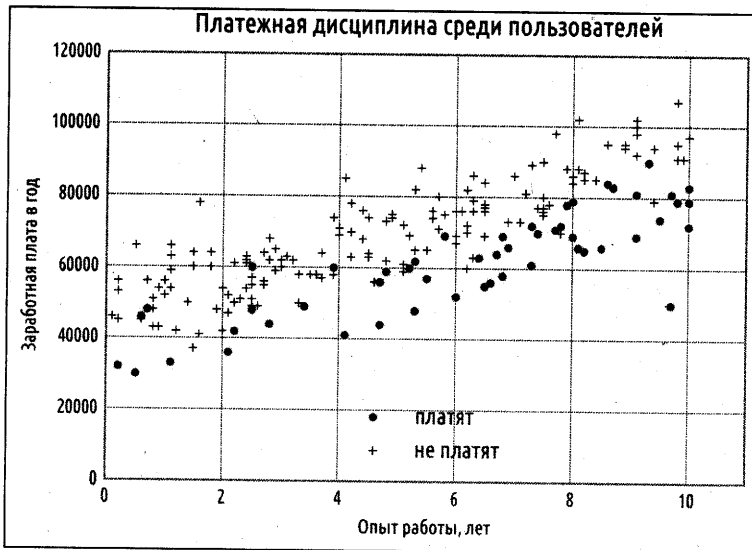


Рис. 16.1. Пользователи, оплатившие и не оплатившие аккаунты

<sup>1</sup> Билл Бэйли (1964) — британский комик, музыкант и актер (см. [https://ru.wikipedia.org/wiki/Бэйли,\\_Билл](https://ru.wikipedia.org/wiki/Бэйли,_Билл)). — Прим. пер.



ные, ей присваивается либо 0 (премиум-аккаунт отсутствует) или 1 (премиум-аккаунт есть).

И как обычно, данные размещены в матрице, где каждая строка — это список [опыт работы, зарплата, аккаунт оплачен?]. Преобразуем ее в формат, который потребуется для решения задачи:

```
x = [[1] + row[:2] for row in data] # каждый элемент = [1, стаж, зарплата]
y = [row[2] for row in data]      # каждый элемент = аккаунт оплачен?
```

Очевидно, сперва следует попытаться применить линейную регрессию и найти наилучшую модель:

$$\text{оплаченный аккаунт} = \beta_0 + \beta_1 \text{стаж} + \beta_2 \text{зарплата} + \epsilon.$$

Естественно, ничто не мешает смоделировать задачу таким образом. Полученные результаты приведены на рис. 16.2:

```
rescaled_x = rescale(x) # применить многомерное шкалирование
beta = estimate_beta(rescaled_x, y) # [0.26, 0.43, -0.43]
predictions = [predict(x_i, beta) for x_i in rescaled_x]
```

```
plt.scatter(predictions, y)
plt.xlabel("прогноз")
plt.ylabel("фактически")
plt.show()
```

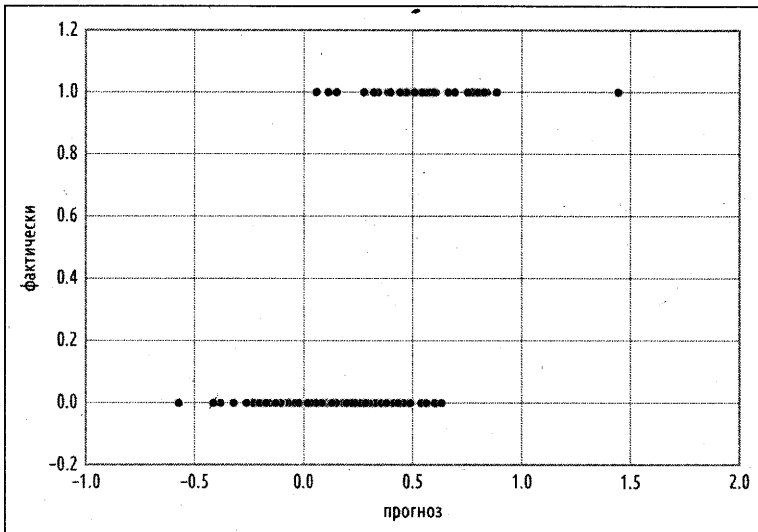


Рис. 16.2. Модель линейной регрессии для предсказания оплаты аккаунтов

Однако этот подход создает пару собственных проблем.

- ◆ Нужно, чтобы получаемые прогнозные значения были 0 или 1, тем самым указывая на принадлежность к классу. Будет замечательно, если они находятся в интервале между 0 и 1, поскольку их можно интерпретировать как вероятно-

сти — значение 0.25 могло бы с вероятностью 25% обозначать оплатившего аккаунт члена соцсети. Однако модель линейной регрессии может порождать результаты в виде огромных положительных или даже отрицательных чисел, которые непонятно как интерпретировать. И действительно, ранее было показано много прогнозных значений с отрицательными значениями.

- ◆ Модель линейной регрессии основана на допущении, что случайные ошибки не коррелируют со столбцами  $x$ . Однако здесь коэффициент регрессии для стажа `experience` равен 0.43, демонстрируя, что увеличение стажа ведет к увеличению правдоподобия привилегированного аккаунта. Это означает, что модель генерирует очень большие значения для лиц с очень большим стажем. Но, как мы знаем, фактические значения должны быть не выше 1, и, значит, неизбежно очень большие результаты (и, следовательно, очень большие значения `experience`) соответствуют очень большим отрицательным значениям случайной ошибки. Поскольку это как раз тот случай, то наша оценка `beta` является смещенной.

Вместо этого нужно, чтобы большие положительные значения `dot(x_i, beta)` соответствовали вероятности, близкой к 1, а большие отрицательные значения — вероятности, близкой к 0. Этого можно добиться, применив к результату еще одну функцию.

## Логистическая функция

В модели логистической регрессии используется *логистическая функция*, чей график показан на рис. 16.3:

```
# логистическая функция
def logistic(x):
    return 1.0 / (1 + math.exp(-x))
```

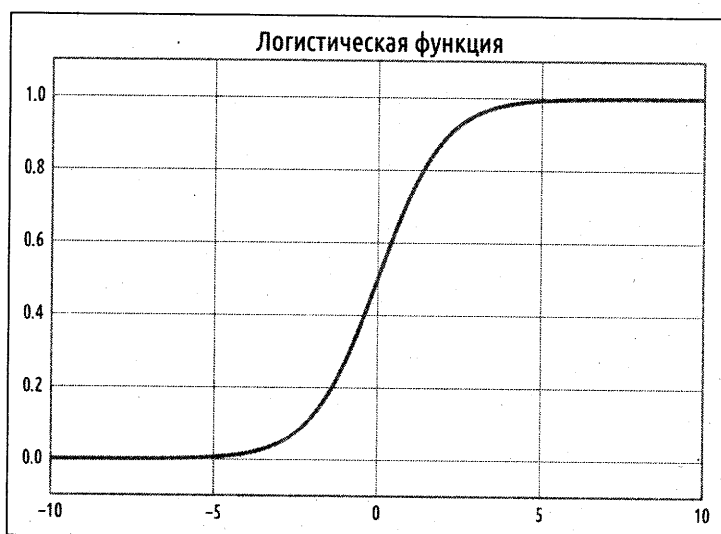


Рис. 16.3. Логистическая функция

Если входящее значение увеличивается и является положительным, то результат функции стремится к 1. Если же входящее значение увеличивается и является отрицательным, то результат стремится к 0 (так называемое сужающее свойство). Дополнительно к этому у нее есть еще одно удобное свойство — ее производная задается в виде:

```
# первая производная логистической функции (штрих)
def logistic_prime(x):
    return logistic(x) * (1 - logistic(x))
```

которой и воспользуемся для подбора модели:

$$y_i = f(x_i, \beta) + \varepsilon_i,$$

где  $f$  — логистическая функция `logistic`.

Напомним, что в случае с линейной регрессией подбор модели осуществлялся путем минимизации суммы квадратов случайных ошибок, после чего был выбран параметр  $\beta$ , который максимизировал правдоподобие данных.

Здесь эти операции не эквивалентны, поэтому воспользуемся методом градиентного спуска, чтобы непосредственно максимизировать правдоподобие. Следовательно, потребуется вычислить функцию правдоподобия и ее градиент.

Согласно модели, при наличии некоторого коэффициента  $\beta$  каждый  $y_i$  должен быть равным 1 с вероятностью  $f(x_i, \beta)$  и 0 с вероятностью  $1 - f(x_i, \beta)$ .

В частности, ДФР для  $y_i$  может быть записана следующим образом:

$$p(y_i | x_i, \beta) = f(x_i, \beta)^{y_i} (1 - f(x_i, \beta))^{1 - y_i},$$

поскольку при  $y_i = 0$  она равна:

$$1 - f(x_i, \beta)$$

а при  $y_i = 1$ :

$$f(x_i, \beta)$$

Оказывается, что на самом деле проще максимизировать *логарифмическую функцию правдоподобия*:

$$\ln L(\beta | x_i, y_i) = y_i \ln f(x_i, \beta) + (1 - y_i) \ln(1 - f(x_i, \beta)).$$

Поскольку логарифм — монотонно возрастающая функция, то любая оценка  $\beta$ , которая максимизирует логарифмическую функцию правдоподобия, также максимизирует саму функцию правдоподобия, и наоборот.

# логарифмическая функция правдоподобия для данных с индексом  $i$

```
def logistic_log_likelihood_i(x_i, y_i, beta):
    if y_i == 1:
        return math.log(logistic(dot(x_i, beta)))
    else:
        return math.log(1 - logistic(dot(x_i, beta)))
```

Если допустить, что разные точки данных друг от друга независимы, то полное правдоподобие — это просто произведение частных правдоподобий, вследствие чего полное логарифмическое правдоподобие — это сумма частных правдоподобий:

```
# логарифмическое правдоподобие
def logistic_log_likelihood(x, y, beta):
    return sum(logistic_log_likelihood_i(x_i, y_i, beta)
               for x_i, y_i in zip(x, y))
```

Несколько правил дифференцирования дадут градиент:

```
# частная по ij
def logistic_log_partial_ij(x_i, y_i, beta, j):
    """здесь i - индекс точки данных, j - индекс производной"""

    return (y_i - logistic(dot(x_i, beta))) * x_i[j]
```

```
# градиент по i
def logistic_log_gradient_i(x_i, y_i, beta):
    """градиент логарифмической функции правдоподобия,
    соответствующий i-й точке данных"""

    return [logistic_log_partial_ij(x_i, y_i, beta, j)
            for j, _ in enumerate(beta)]
```

```
# градиент
def logistic_log_gradient(x, y, beta):
    return reduce(vector_add,
                  [logistic_log_gradient_i(x_i, y_i, beta)
                   for x_i, y_i in zip(x, y)])
```

На данный момент имеются все необходимые составляющие.

## Применение модели

Сначала необходимо разделить данные на обучающую и контрольную выборки:

```
random.seed(0)
x_train, x_test, y_train, y_test = train_test_split(rescaled_x, y, 0.33)
```

```
# требуется максимизировать функцию логарифмического правдоподобия
# на обучающих данных
fn = partial(logistic_log_likelihood, x_train, y_train)
gradient_fn = partial(logistic_log_gradient, x_train, y_train)
```

```
# выбрать произвольную отправную точку
beta_0 = [random.random() for _ in range(3)]
```

```
# и максимизировать методом градиентного спуска
beta_hat = maximize_batch(fn, gradient_fn, beta_0)
```

С другой стороны, можно применить стохастический градиентный спуск:

```
# бета "с крышкой" означает оценку параметра бета
beta_hat = maximize_stochastic(logistic_log_likelihood_i,
                               logistic_log_gradient_i,
                               x_train, y_train, beta_0)
```

В обоих случаях приближенно находим:

```
beta_hat = [-1.90, 4.05, -3.87]
```

Они представляют собой коэффициенты для шкалированных данных `rescaled`, которые, тем не менее, можно привести обратно к исходному виду:

```
beta_hat_unscaled = [7.61, 1.42, -0.000249]
```

К сожалению, эти данные сложнее интерпретировать по сравнению с коэффициентами линейной регрессии. При прочих равных дополнительный год стажа работы добавляет 1.42 к входящему аргументу в логистическую функцию `logistic`, а дополнительные \$10 000 в зарплате при прочих равных отнимает 2.49 из входящего аргумента.

Однако на результат влияют и другие входящие значения. Если значение `dot(beta, x_i)` уже большое (соответствуя вероятности, близкой к 1), то даже его значительное увеличение не сможет сильно повлиять на вероятность. Если же это число близко к 0, то даже незначительное его увеличение может существенно увеличить вероятность<sup>2</sup>.

Можно только сказать, что — при прочих равных условиях — чем больше стаж работы, тем больше вероятность, что аккаунты будут оплачиваться, и что — при прочих равных условиях — чем выше зарплата, тем меньше вероятность, что они будут оплачиваться. (Это также было очевидно после построения графика.)

## Качество подбора модели

Теперь воспользуемся контрольной выборкой. Посмотрим, что получится, если предсказывать *оплату аккаунта* всякий раз, когда вероятность превышает 0.5:

```
# переменные для подсчета соответственно истинноположительных,
# ложноположительных, истинноотрицательных
# и ложноотрицательных результатов
true_positives = false_positives = true_negatives = false_negatives = 0

for x_i, y_i in zip(x_test, y_test):
    predict = logistic(dot(beta_hat, x_i))
```

<sup>2</sup> Для малых значений логистическая кривая имеет сильный наклон, дающий большое усиление. Когда значение становится больше, усиление падает. — *Прим. пер.*

```

if y_i == 1 and predict >= 0.5: # TP: оплачено и прогноз тот же
    true_positives += 1
elif y_i == 1: # FN: оплачено, но прогноз обратный
    false_negatives += 1
elif predict >= 0.5: # FP: неоплачено, но прогноз обратный
    false_positives += 1
else: # TN: неоплачено и прогноз тот же
    true_negatives += 1

```

```
# точность и полнота
```

```
precision = true_positives / (true_positives + false_positives)
```

```
recall = true_positives / (true_positives + false_negatives)
```

В результате будет получен достаточно солидный результат: точность 93% ("прогноз об *оплате аккаунта* правдив в 93% случаев") и полнота 82% ("прогноз об *оплате аккаунта* делается в отношении 82% пользователей, оплативших аккаунт"), что в обоих случаях является заслуживающим уважения результатом.

Кроме того, можно сопоставить прогнозы и фактические данные на графике (рис. 16.4), который также иллюстрирует хорошую производительность модели:

```
# предсказания
```

```
predictions = [logistic(dot(beta_hat, x_i)) for x_i in x_test]
```

```
plt.scatter(predictions, y_test)
```

```
plt.xlabel("Прогнозная вероятность")
```

```
plt.ylabel("Фактический результат")
```

```
plt.title("Прогноз и фактические данные")
```

```
plt.show()
```

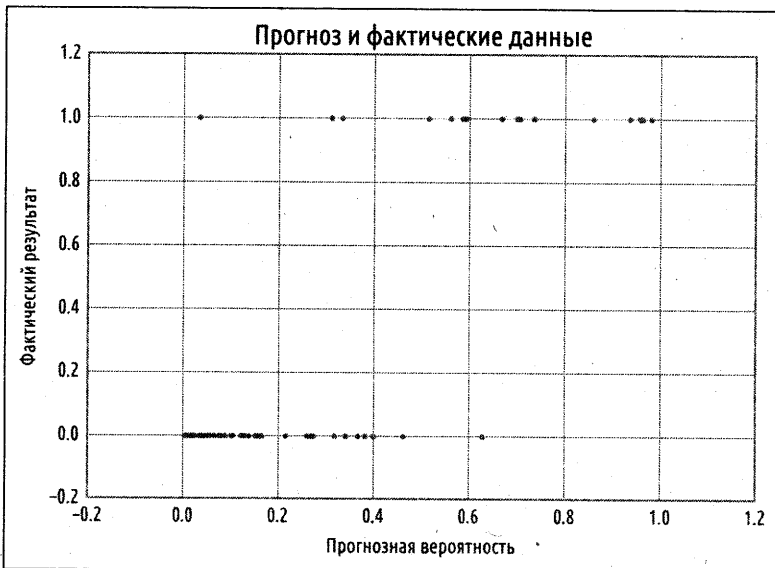


Рис. 16.4. Прогноз на основе логистической регрессии в сравнении с фактическими данными

## Метод опорных векторов

Набор точек, где  $\text{dot}(\beta_{\text{hat}}, x_i) = 0$ , представляет границу между классами. Можно построить диаграмму, чтобы увидеть, что конкретно модель делает (рис. 16.5).

Эта граница является *гиперплоскостью*, которая разделяет пространство параметров на два полупространства, соответствующие *прогнозу оплаты* и *прогнозу неоплаты*. Это было обнаружено в виде побочного эффекта обнаружения наиболее правдоподобной логистической модели.

Альтернативный подход к классификации заключается в поиске гиперплоскости, которая "наилучшим образом" разделяет классы в обучающих данных. Эта идея лежит в основе *метода опорных векторов* (support vector machine), который находит такую гиперплоскость, которая максимизирует расстояние до ближайшей точки в каждом классе (рис. 16.6).

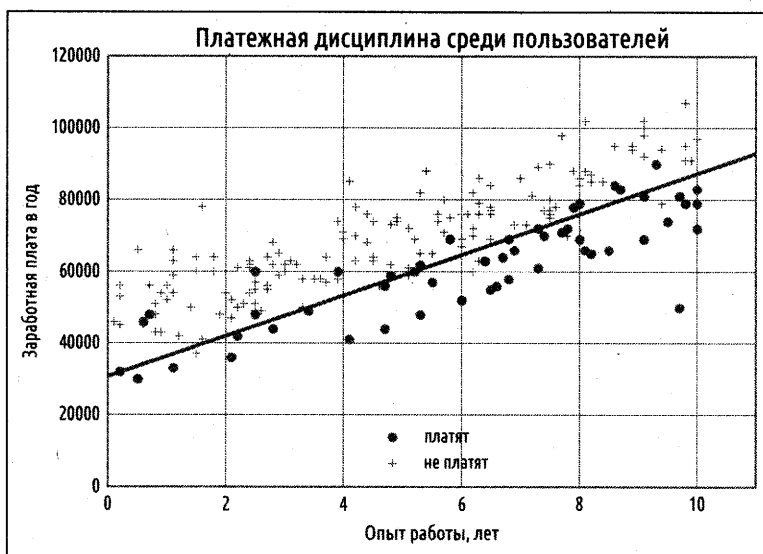


Рис. 16.5. Оплатившие и неоплатившие пользователи с решающей границей

Нахождение такой гиперплоскости является задачей оптимизации, которая подразумевает использование приемов, находящихся на более продвинутом уровне. Другая проблема заключается в том, что разделяющая гиперплоскость может не существовать вообще. В наборе данных об оплате аккаунтов просто отсутствует линия, которая идеально отделяет оплативших пользователей от неоплативших.

Эту проблему можно (иногда) обойти путем перевода данных в пространство с более высокой размерностью. Например, рассмотрим простой одномерный набор данных, показанный на рис. 16.7.

Ясно, что гиперплоскость, которая могла бы отделять положительные образцы от отрицательных, отсутствует. Однако посмотрим, что произойдет, если отобразить



Рис. 16.6. Разделяющая гиперплоскость

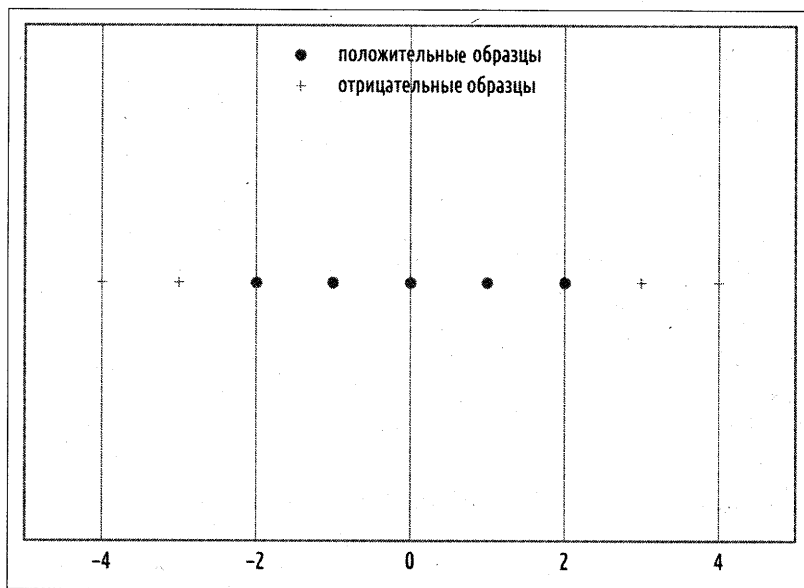


Рис. 16.7. Неразделимый одномерный набор данных

этот набор данных в две размерности, направив точку  $x$  в  $(x, x^2)$ . Неожиданно обнаружится, что можно найти гиперплоскость, разделяющую данные (рис. 16.8).

Обычно этот прием называют *ядерным трюком* (kernel trick), т. к. вместо фактического отображения точек в пространство более высокой размерности (что может быть затратным, если точек много, и отображение осложнено) используется "ядерная" функция, которая вычисляет скалярные произведения в пространстве большей размерности и использует их результаты для поиска гиперплоскости.



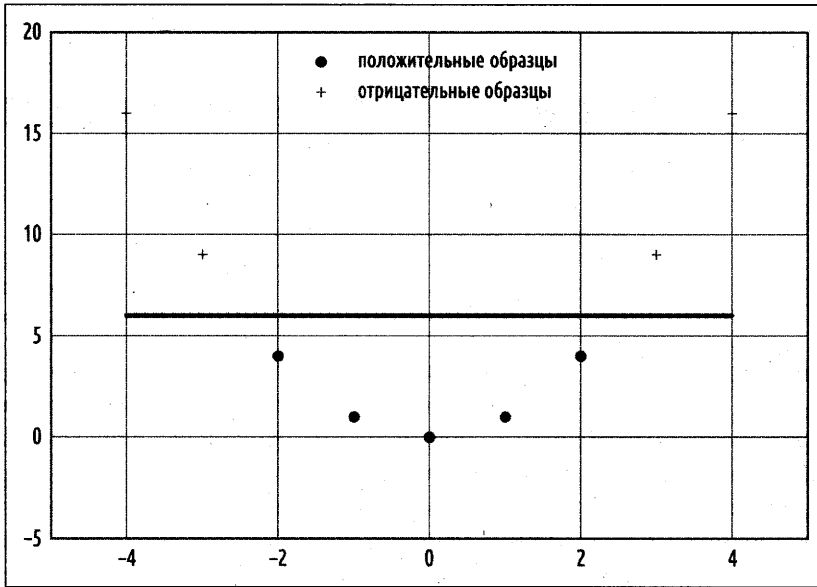


Рис. 16.8. Набор данных становится разделимым в более высокой размерности

Трудно (и, вероятно, неразумно) использовать метод опорных векторов без специализированного оптимизационного программного обеспечения, написанного специалистами с соответствующим опытом, поэтому здесь на изучении данной темы придется поставить точку.

## Для дальнейшего изучения

- ◆ Библиотека `scikit-learn` располагает модулями для решения задач на основе логистической регрессии<sup>3</sup> и метода опорных векторов<sup>4</sup>.
- ◆ Библиотека `libsvm` (<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>) реализует метод опорных векторов, который лежит в основе его реализации в `scikit-learn`. Веб-сайт библиотеки располагает большим количеством разнообразной полезной документации, касающейся этого метода.

<sup>3</sup> [http://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](http://scikit-learn.org/stable/modules/linear_model.html#logistic-regression).

<sup>4</sup> <http://scikit-learn.org/stable/modules/svm.html>.

# Деревья принятия решений

Дерево — это непостижимая загадка.  
Джим Вудринг<sup>1</sup>

Директор подразделения по поиску талантов с переменным успехом провел ряд собеседований с претендентами, чьи заявления поступили с сайта DataSciencester, и в результате собрал набор данных, состоящий из нескольких (качественных) признаков, описывающих каждого претендента, а также о том, прошел ли претендент собеседование успешно или нет. Он интересуется, можно ли использовать эти данные для построения модели, определяющей, какие претенденты проходят собеседование успешно, благодаря чему ему не придется тратить время на проведение собеседований.

Эта задача, оказывается, хорошо подходит для модели на основе еще одного инструмента прогнозного моделирования в наборе инструментов аналитика данных, именуемого *деревьями принятия решений* (ДПР, decision tree) или *решающими деревьями*.

## Что такое дерево принятия решений?

В модели деревьев принятия решений используется древовидная структура данных, которая предоставляет несколько возможных *путей принятия решения* и исход для каждого пути.

Тот, кто когда-либо играл в игру "20 вопросов"<sup>2</sup>, уже знаком с деревьями принятия решений, даже не подозревая об этом. Вот пример.

— Я думаю о животном.

— У него больше пяти лап?

— Нет.

— Он вкусный?

— Нет.

— Его можно увидеть на обратной стороне австралийского 5-центовика?

— Да.

<sup>1</sup> Джим Вудринг (1952) — американский мультипликатор, художник изящных искусств, писатель и дизайнер игрушек (см. [https://en.wikipedia.org/wiki/Jim\\_Woodring](https://en.wikipedia.org/wiki/Jim_Woodring)). — Прим. пер.

<sup>2</sup> [https://en.wikipedia.org/wiki/Twenty\\_Questions](https://en.wikipedia.org/wiki/Twenty_Questions).

— Это ехидна?

— Да, ты угадал!

Это соответствует пути:

"Не более 5 лап" → "Не вкусный" → "На 5-центовой монете" → "Ехидна!"  
 на своеобразном (и не очень подробном) дереве принятия решений "Угадать животное" (рис. 17.1).

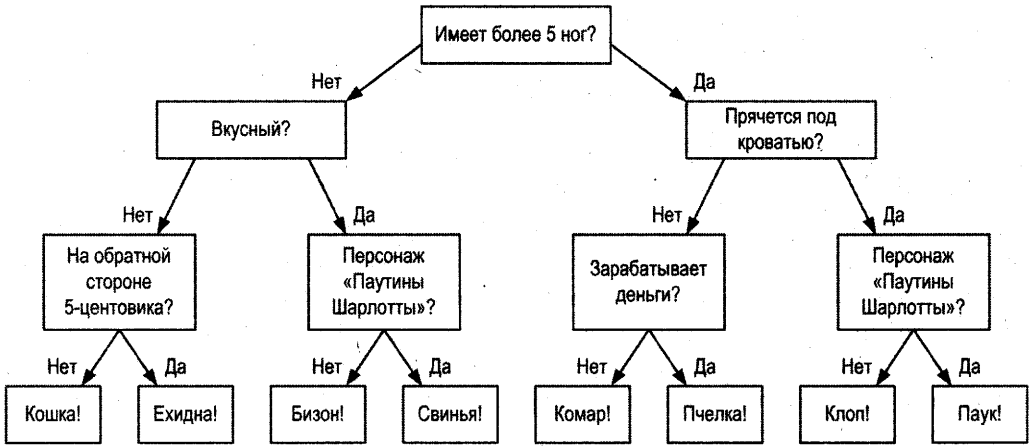


Рис. 17.1. Дерево принятия решений "Угадать животное"

Модели на основе деревьев принятия решений обладают многими преимуществами. Они понятны и легко интерпретируемы, а процесс, на основе которого они достигают прогноза, абсолютно прозрачен. В отличие от других моделей, которые были рассмотрены до сих пор, модели на основе деревьев принятия решений легко справляются с сочетанием численных (количество ног) и категориальных (вкусный/не вкусный) признаков (или атрибутов, что то же самое) и даже могут классифицировать данные, для которых атрибуты отсутствуют.

Вместе с тем нахождение "оптимального" дерева принятия решений для набора обучающих данных — задача вычислительно очень сложная. (Эта проблема будет обойдена путем построения приемлемого дерева вместо оптимального, хотя при наличии больших наборов данных его создание все равно может потребовать значительной обработки.) Важнее все же то, что очень легко (и это очень плохо) построить деревья принятия решений, которые *избыточно подогнаны* под обучающую выборку и которые недостаточно хорошо обобщают не встречавшиеся ранее данные. Мы обратимся к способам решения этих проблем.

Обычно деревья принятия решений подразделяются на *классификационные деревья* (которые генерируют категориальные значения) и *регрессионные деревья* (которые генерируют численные или непрерывные значения). В этой главе мы сосредоточимся на классификационных деревьях и коснемся алгоритма ID3 для вычисления дерева принятия решений в результате обучения набором маркированных данных, который поможет понять на практике принцип работы ДПР. Чтобы упростить ра-

боту, ограничимся задачами с бинарными результатами типа "Следует ли нанять этого претендента?" или "Какую из двух реклам следует показать этому посетителю сайта: А или В?", или "Будет ли меня тошнить, если съесть эту пищу из офисного холодильника?"

## Энтропия

Для того чтобы построить дерево принятия решений, необходимо решить, какие вопросы задавать и в каком порядке. На каждом уровне дерева имеются некоторые возможности, которые были исключены и которые еще остались. Получив информацию о том, что у животного не более пяти ног, исключают возможность, что это кузнечик, но оставлен вариант с уткой. Каждый возможный вопрос разделяет оставшиеся возможные варианты в соответствии с полученным ответом.

В идеале хотелось бы подобрать вопросы, ответы на которые дают максимум информации о предмете предсказания. Если есть общий вопрос ("да/нет"), для которого ответы "да" всегда соответствуют истинным результатам (True), а ответы "нет" — ложным (False), или наоборот, то это потрясающий вопрос для ДПР. И с другой стороны, общий вопрос, для которого ни один из ответов не предоставляет достаточного количества новой информации о предмете предсказания, вероятно, был подобран не очень хорошо.

Идея о "количестве информации" объясняется термином *энтропия*. Это слово часто приводят для обозначения беспорядка, хаоса. Здесь оно используется для обозначения неопределенности или непредсказуемости, связанной с данными.

Пусть имеется множество  $S$  данных, каждый маркированный элемент которого принадлежит одному из классов конечного множества классов  $C_1, \dots, C_n$ . Если все точки данных принадлежат только одному классу, то неопределенность фактически отсутствует, вследствие чего требуется низкая энтропия. Если же точки данных распределены по классам равномерно, то существует большая неопределенность, и значит, требуется высокая энтропия.

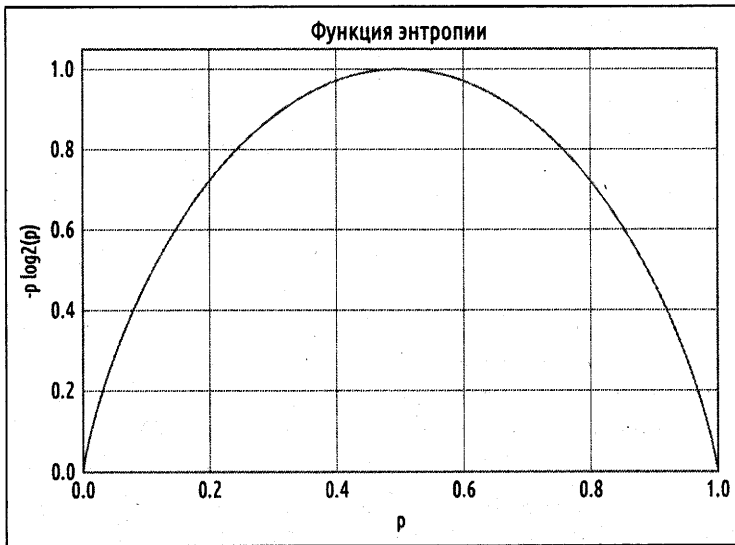
На языке математики, если  $p_i$  — это пропорция данных, помеченная как класс  $c_i$ , то энтропия определяется следующей формулой<sup>3</sup>:

$$H(S) = -p_1 \log_2 p_1 - \dots - p_n \log_2 p_n$$

со (стандартным) соглашением, что  $0 \cdot \log_2 0 = 0$ .

Особо не вдаваясь в жуткие подробности, лишь отметим, что каждый член  $-p_i \log_2 p_i$  принимает неотрицательное значение, близкое к нулю, именно тогда, когда параметр  $p_i$  близок к нулю или единице (рис. 17.2).

<sup>3</sup> Речь идет о вероятностной мере неопределенности Шеннона (или информационной энтропии). Данная формула означает, что прирост информации равен утраченной неопределенности. — *Прим. пер.*

Рис. 17.2. График кривой  $-p \log_2 p$ 

Следовательно, энтропия будет низкой, когда каждый  $p_i$  близок к 0 или к 1 (т. е. когда подавляющая часть данных находится в одном классе), и напротив, она будет выше при наличии нескольких  $p_i$ , чьи значения не близки к 0 (т. е. когда данные распределены между несколькими классами). Это именно то поведение, которое требуется.

Все это достаточно легко реализовать в одной функции:

```
def entropy(class_probabilities):
    """при заданном списке вероятностей классов вычислить энтропию"""
    return sum(-p * math.log(p, 2)
              for p in class_probabilities
              if p)           # игнорировать нулевые вероятности
```

Данные состоят из пар (вход, метка), а значит, придется вычислить вероятности классов непосредственно. Заметим, что на самом деле интерес вызывает не то, какая метка связана с конкретной вероятностью, а только то, каковы вероятности:

```
# вероятности классов
def class_probabilities(labels):
    total_count = len(labels)           # суммарное число
    return [count / total_count
            for count in Counter(labels).values()]
```

```
# энтропия маркированных данных
def data_entropy(labeled_data):
    labels = [label for _, label in labeled_data]
    probabilities = class_probabilities(labels)
    return entropy(probabilities)
```

## Энтропия разбиения

Ранее в этой главе была вычислена энтропия (читай "неопределенность") одиночного набора маркированных данных. Далее, каждый шаг дерева принятия решений предполагает постановку вопроса, ответ на который расщепляет данные на один или (надо полагать) более подмножеств. Например, вопрос "у него больше пяти лап?" разделяет животных на тех, у которых их больше пяти (пауки) и остальных (ехидна).

Следовательно, требуется некое представление об энтропии, получаемой в результате разбиения набора данных определенным образом. Энтропия разбиения должна быть низкой, если данные расщепляются на подгруппы, которые сами обладают низкой энтропией (т. е. очень определенные), и высокой, если внутри имеются подгруппы (которые являются большими и) с высокой энтропией (т. е. крайне неопределенные).

Например, глупый, но довольно удачный вопрос об австралийской 5-центовой монете разделил оставшихся на тот момент животных на подгруппы  $S_1 = \{\text{ехидна}\}$  и  $S_2 = \{\text{все остальные}\}$ , где  $S_2$  — большая подгруппа с высокой энтропией. (У  $S_1$  энтропия отсутствует, впрочем оно представляет лишь малую часть оставшихся "классов.")

С математической точки зрения, если расщеплять данные  $S$  на подгруппы  $S_1, \dots, S_m$ , содержащие пропорции  $q_1, \dots, q_m$  данных, то энтропия разбиения вычисляется в виде взвешенной суммы:

$$H = q_1 H(S_1) + \dots + q_m H(S_m),$$

что можно реализовать следующим образом:

```
# энтропия разбиения для подмножеств
def partition_entropy(subsets):
    """найти энтропию исходя из этого разбиения данных на подгруппы
    подгруппы представляют собой список списков маркированных данных"""

    total_count = sum(len(subset) for subset in subsets)

    return sum( data_entropy(subset) * len(subset) / total_count
                for subset in subsets )
```



Проблема с этим подходом заключается в том, что расщепление по атрибуту с несколькими разными значениями приведет к очень низкой энтропии, из-за переобучения. Например, представим, что вы работаете в банке и пытаетесь построить дерево принятия решений, чтобы предсказать, какие из клиентов банка, скорее всего, перестанут платить по ипотеке, используя для этого в качестве обучающей выборки некие исторические данные. Допустим далее, что набор данных содержит номер социального страхования (НСС) каждого клиента. Расщепление по НСС будет генерировать подмножества, состоящие из одного человека, и каждое неизбежно будет иметь нулевую энтропию. Однако модель, которая опирается на НСС, определенно не сможет обобщать за пределами обучающей выборки. По этой причине при создании ДПР, очевидно, следует пытаться избегать (или, когда нужно, разбивать на интервалы) атрибуты с большим числом возможных значений.

## Создание дерева принятия решений

Директор предоставил данные о прошедших собеседовании претендентах, которые состоят (согласно вашей спецификации) из пар (input, label), где input — это словарь dict с атрибутами претендента и label — метка со значением True (претендент прошел собеседование успешно) либо False (претендент собеседование не прошел). В качестве атрибутов представлены уровень претендента, наиболее предпочитаемый язык программирования, общается ли он в Twitter и имеет ли он ученую степень.

```
inputs = [
    ({'level':'Senior', 'lang':'Java', 'tweets':'no', 'phd':'no'}, False),
    ({'level':'Senior', 'lang':'Java', 'tweets':'no', 'phd':'yes'}, False),
    ({'level':'Mid', 'lang':'Python', 'tweets':'no', 'phd':'no'}, True),
    ({'level':'Junior', 'lang':'Python', 'tweets':'no', 'phd':'no'}, True),
    ({'level':'Junior', 'lang':'R', 'tweets':'yes', 'phd':'no'}, True),
    ({'level':'Junior', 'lang':'R', 'tweets':'yes', 'phd':'yes'}, False),
    ({'level':'Mid', 'lang':'R', 'tweets':'yes', 'phd':'yes'}, True),
    ({'level':'Senior', 'lang':'Python', 'tweets':'no', 'phd':'no'}, False),
    ({'level':'Senior', 'lang':'R', 'tweets':'yes', 'phd':'no'}, True),
    ({'level':'Junior', 'lang':'Python', 'tweets':'yes', 'phd':'no'}, True),
    ({'level':'Senior', 'lang':'Python', 'tweets':'yes', 'phd':'yes'}, True),
    ({'level':'Mid', 'lang':'Python', 'tweets':'no', 'phd':'yes'}, True),
    ({'level':'Mid', 'lang':'Java', 'tweets':'yes', 'phd':'no'}, True),
    ({'level':'Junior', 'lang':'Python', 'tweets':'no', 'phd':'yes'}, False)
]
```

Дерево будет состоять из *решающих узлов* (которые содержат вопрос и направляют в разные стороны в зависимости от ответа) и *листовых узлов* (которые делают прогноз). Дерево будет построено при помощи сравнительно простого *алгоритма ID3*, который работает следующим образом. Допустим, заданы некоторые маркированные данные и список атрибутов, которые необходимо учитывать при принятии решения о ветвлении.

1. Если все данные имеют одинаковую метку, то создать листовой узел, который предсказывают эту метку, и остановиться.
2. Если список атрибутов пуст (т. е. больше вопросов не осталось), то создать листовой узел, который предсказывает наиболее общую метку, и затем остановиться.
3. В противном случае попытаться разбить данные по каждому атрибуту.
4. Выбрать разбиение с наименьшей энтропией.
5. Добавить решающий узел на основе выбранного атрибута.
6. Рекурсивно повторить для каждой подгруппы, получившейся в результате разбиения, используя оставшиеся атрибуты.

Такой алгоритм называется "жадным", поскольку на каждом шаге он непосредственно выбирает наиболее оптимальный вариант. Однако он не учитывает случай,

когда для заданного набора данных есть более оптимальное дерево, но с менее оптимальным первым шагом. В этом случае алгоритм его не найдет. Тем не менее, алгоритм сравнительно понятен и легко осуществим, что делает его хорошей отправной точкой для исследования деревьев принятия решений.

Пройдем эти шаги "в ручном режиме", используя набор данных о претендентах, которые прошли собеседование. Набор данных содержит метки со значениями True и False и четыре атрибута, по которым может выполняться расщепление. Итак, первый шаг будет заключаться в поиске разбиения с наименьшей энтропией. Начнем с написания функции, которая выполняет разбиение:

```
# разбить входящие данные по атрибуту
def partition_by(inputs, attribute):
    """каждый элемент в inputs представляет собой пару
    (словарь attribute_dict, метка label).
    Возвращает словарь dict,
    где ключ = значение атрибута attribute_value, значение = inputs"""
    groups = defaultdict(list)
    for input in inputs:
        key = input[0][attribute] # взять значение атрибута и добавить
        groups[key].append(input) # входящую пару к соответствующему списку
    return groups
```

и функции, которая использует ее для вычисления энтропии:

```
# энтропия разбиения по атрибуту
def partition_entropy_by(inputs, attribute):
    """вычисляет энтропию, соответствующую заданному разбиению"""
    partitions = partition_by(inputs, attribute)
    return partition_entropy(partitions.values())
```

Затем нужно найти разбиение с минимальной энтропией относительно всего набора данных:

```
for key in ['level', 'lang', 'tweets', 'phd']:
    print(key, partition_entropy_by(inputs, key))
```

```
# level 0.693536138896
# lang 0.860131712855
# tweets 0.788450457308
# phd 0.892158928262
```

Минимальная энтропия получается из расщепления по атрибуту level (уровень), после чего для каждого возможного значения атрибута level нужно создать поддерево. Каждый претендент со значением Mid (средний) этого атрибута помечен как True, и, следовательно, поддерево Mid — это всего лишь листовая узел, который предсказывает True. У претендентов со значением Senior (высокий) этого атрибута имеется смесь значений True и False, поэтому следует снова выполнить расщепление:



```

senior_inputs = [(input, label)
                 for input, label in inputs if input["level"] == "Senior"]

for key in ['lang', 'tweets', 'phd']:
    print(key, partition_entropy_by(senior_inputs, key))

# lang 0.4
# tweets 0.0
# phd 0.950977500433
    
```

Результат показывает, что следующее расщепление должно пройти по атрибуту tweets, которое дает разбиение с нулевой энтропией. Для претендентов с уровнем Senior атрибут tweets со значением "да" всегда дает True, а со значением "нет" — всегда False.

Наконец, если проделать то же самое для претендентов с уровнем Junior (начальный), то в конечном итоге разделение пройдет по атрибуту phd (ученая степень), после чего обнаружится, что отсутствие степени всегда дает True и наличие степени — всегда False.

На рис. 17.3 показано полное дерево принятия решений.

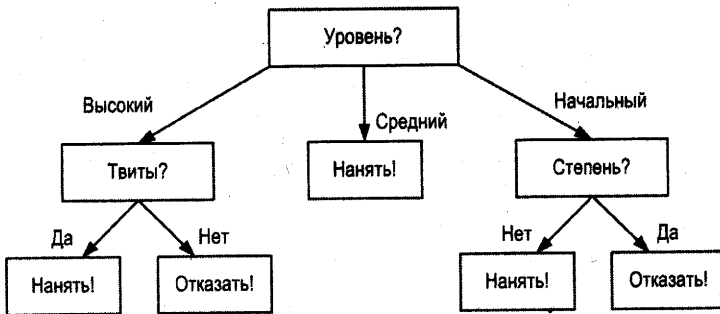


Рис. 17.3. Дерево принятия решений для найма сотрудников

## Собираем все вместе

После того как работа алгоритма продемонстрирована, следует реализовать его в более общем плане. Иными словами, следует решить, каким образом представлять деревья. Воспользуемся наиболее легким представлением из возможных. Определим *дерево* как одну из следующих альтернатив:

- ◆ True;
- ◆ False;
- ◆ кортеж (attribute, subtree\_dict),

где True — листовой узел, который для любого входящего значения возвращает истину; False — листовой узел, который для любого входящего значения возвращает ложь; кортеж — решающий узел, который для любого входящего значения находит

его значение атрибута `attribute` и классифицирует входящее значение, используя соответствующее поддерево.

В таком представлении дерево найма сотрудников примет следующий вид:

```
('level',
 {'Junior': ('phd', {'no': True, 'yes': False}),
  'Mid': True,
  'Senior': ('tweets', {'no': False, 'yes': True})})
```

Остался вопрос: что делать при встрече неожиданного (или отсутствующего) значения атрибута? Что дерево найма сотрудников должно делать, если встретится претендент, чей атрибут `level` равен "Intern" (стажер)? Этот случай обрабатывается добавлением ключа `None`, который всего лишь предсказывает наиболее общую метку. (Правда, так не стоит делать в случаях, когда `None` — это значение, которое имеется в самих данных.)

При таком представлении входящие значения классифицируются следующим образом:

```
def classify(tree, input):
    """классифицировать входящие значения, используя заданное дерево ДДР"""

    # если это листовой узел, вернуть его значение
    if tree in [True, False]:
        return tree

    # иначе дерево состоит из атрибута, по которому пройдет расщепление,
    # и словаря, где ключи - это значения этого атрибута, а
    # значения - это поддерева, подлежащие рассмотрению в следующую очередь
    attribute, subtree_dict = tree

    subtree_key = input.get(attribute) # None, если на входе
                                        # отсутствующий атрибут

    if subtree_key not in subtree_dict: # если для ключа нет поддерева,
        subtree_key = None              # использовать поддерево None

    subtree = subtree_dict[subtree_key] # выбрать соответствующее поддерево
    return classify(subtree, input)     # и использовать для классификации
```

Осталось только построить из обучающей выборки древовидное представление:

```
# построить дерево на основе алгоритма ID3
def build_tree_id3(inputs, split_candidates=None):

    # если это первый проход, то
    # все ключи первоначальных входящих данных - это выделенные претенденты
    if split_candidates is None:
        split_candidates = inputs[0][0].keys()
```

```

# подсчитать число True и False во входящих значениях
num_inputs = len(inputs)
num_trues = len([label for item, label in inputs if label])
num_falses = num_inputs - num_trues

if num_trues == 0: return False # если True отсутствуют,
                               # вернуть лист False
if num_falses == 0: return True # если False отсутствуют,
                                # вернуть лист True

if not split_candidates:      # если больше нет кандидатов
    return num_trues >= num_falses # вернуть лист большинства

# в противном случае выполнить расщепление по лучшему атрибуту
best_attribute = min(split_candidates,
                     key=partial(partition_entropy_by, inputs))

partitions = partition_by(inputs, best_attribute)
new_candidates = [a for a in split_candidates
                  if a != best_attribute]

# рекурсивно создать поддерева
subtrees = { attribute_value : build_tree_id3(subset, new_candidates)
            for attribute_value, subset in partitions.iteritems() }

subtrees[None] = num_trues > num_falses # случай по умолчанию

return (best_attribute, subtrees)

```

Каждый лист в построенном дереве состоит целиком из входящих значений True или False. Следовательно, дерево прекрасно справляется с предсказанием на основе обучающей выборки. Кроме того, его можно применить к новым данным, которые изначально в обучающей выборке отсутствовали:

```

tree = build_tree_id3(inputs)

classify(tree, { "level" : "Junior",
                 "lang" : "Java",
                 "tweets" : "yes",
                 "phd" : "no" } ) # True

classify(tree, { "level" : "Junior",
                 "lang" : "Java",
                 "tweets" : "yes",
                 "phd" : "yes" } ) # False

```

А также к данным с пропущенными или неожиданными значениями:

```

classify(tree, { "level" : "Intern" } ) # True
classify(tree, { "level" : "Senior" } ) # False

```



Поскольку задача в основном, заключается в том, чтобы продемонстрировать процесс создания дерева, то при его построении использовался весь набор данных целиком. Как и прежде, если бы стояла задача создать реальную модель, то (было бы собрано больше данных и) данные были бы разделены на обучающее, проверочное и контрольное подмножества.

## Случайные леса

Поскольку деревья принятия решений могут приспособливаться к обучающим данным, неудивительно, что они имеют тенденцию к переобучению. Один из способов, который помогает избежать этого, — статистический метод *случайных лесов* (random forests). Данный метод подразумевает создание нескольких ДДР, которые затем принимают коллективное решение на основе голосования о том, каким образом классифицировать входящие значения:

```
# классифицировать на основе леса
def forest_classify(trees, input):
    votes = [classify(tree, input) for tree in trees]
    vote_counts = Counter(votes)
    return vote_counts.most_common(1)[0][0]
```

Процесс построения дерева имеет детерминированный характер. Тогда каким образом можно получить случайные деревья?

Прежде всего для этих целей к данным применяются бутстрап-методы (см. разд. "Отступление: бутстрапирование данных" главы 15). Вместо того чтобы обучать деревья на всех элементах `inputs` обучающей выборки, каждое дерево обучают на результате, получаемом после вызова `bootstrap_sample(inputs)`. Поскольку теперь каждое дерево строится на разных данных, деревья будут друг от друга отличаться. (Дополнительное преимущество заключается в том, что для проверки каждого дерева совершенно справедливо можно использовать невыборочные данные, а значит, в качестве обучающей выборки можно смело воспользоваться всеми данными целиком, если разумно подойти к измерению производительности.) Этот метод известен как *бутстрап-агрегирование* или *бэггинг*<sup>4</sup>.

Помимо этого, еще одним источником случайности является изменение способа выбора наилучшего атрибута `best_attribute`, по которому выполняется расщепление. Вместо исчерпывающего просмотра всех оставшихся атрибутов сначала отбирается их случайное подмножество, а затем расщепление идет по наилучшему из них:

```
# если уже осталось мало претендентов, обратиться ко всем
if len(split_candidates) <= self.num_split_candidates:
    sampled_split_candidates = split_candidates
```

<sup>4</sup> Бутстрап-агрегирование или бэггинг — метод формирования ансамблей классификаторов с использованием случайной выборки с возвратом. При формировании выборок из исходного набора данных случайным образом отбирается несколько подмножеств такого же размера, что и исходный. Затем на основе каждой строится классификатор и их выходы комбинируются путем голосования или простого усреднения. — *Прим. пер.*

```
# в противном случае сделать случайную выборку
else:
    sampled_split_candidates = random.sample(split_candidates,
                                             self.num_split_candidates)

# теперь выбрать наилучший атрибут только из этих претендентов
best_attribute = min(sampled_split_candidates,
                     key=partial(partition_entropy_by, inputs))

partitions = partition_by(inputs, best_attribute)
```

Это пример более широкой методики, называемой *ансамблевым обучением*, когда в целях получения единой усиленной модели объединяются несколько *слабых* моделей, дающих низкие результаты обучения (как правило, это модели с высоким смещением и низкой дисперсией).

Модели на основе случайных лесов являются одними из самых популярных и универсальных.

## Для дальнейшего изучения

- ◆ Библиотека `scikit-learn` располагает множеством моделей на основе дерева принятия решений (<http://scikit-learn.org/stable/modules/tree.html>). Помимо этого, в библиотеке есть модуль `ensemble` (<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.ensemble>), в который включен классификатор на основе решающих деревьев `RandomForestClassifier` и другие ансамблевые методы.
- ◆ В данной главе тема деревьев принятия решений и связанных с ними алгоритмов была затронута лишь вскользь. Википедия является хорошей отправной точкой для более широкого исследования данной темы ([https://en.wikipedia.org/wiki/Decision\\_tree\\_learning](https://en.wikipedia.org/wiki/Decision_tree_learning)).

# Нейронные сети

Люблю бессмыслицы, они пробуждают мозговые клетки.  
*Доктор Сьюз<sup>1</sup>*

*Искусственная нейронная сеть* (artificial neural network) или просто нейросеть — это прогнозная модель, основанием для разработки которой послужил способ организации и принцип функционирования головного мозга. Представим его в виде совокупности нейронов, которые соединены между собой. Каждый нейрон просматривает сигналы, выходящие из других соединенных с ним нейронов, взвешивает их и затем либо реагирует (если результат взвешивания превышает некоторое пороговое значение), либо нет (если не превышает).

Искусственные нейронные сети соответственно состоят из искусственных нейронов, которые производят аналогичное взвешивание своих входящих значений. Нейронные сети могут решать широкий спектр задач, среди которых такие как распознавание рукописного текста и распознавание лиц. Помимо этого, они интенсивно применяются в одной из самых сверхсовременных отраслей науки о данных — глубоком обучении. Однако большинство нейронных сетей представляют собой "черные ящики" — изучение подробностей их устройства не особо прибавит к пониманию того, *каким образом* они решают задачу. К тому же большие нейронные сети могут быть труднообучаемыми. Для большинства задач, с которыми вы будете сталкиваться в качестве начинающего аналитика данных, они, вероятно, едва подойдут. Однако в один прекрасный день, когда вы попытаетесь построить искусственный интеллект, чтобы осуществить компьютерную сингулярность<sup>2</sup>, они окажутся как нельзя кстати.

## Перцептроны

Простейшая модель нейронной сети — это однослойный *перцептрон*, который представляет собой приближенную математическую модель одиночного нейрона, состоящего из  $n$  бинарных входящих сигналов или значений. Он вычисляет взвешенную сумму входящих сигналов и "активизируется", если эта сумма больше или равна нулю:

---

<sup>1</sup> Теодор Сьюз Гейзель (1904–1991) — американский детский писатель и мультипликатор. Автор таких забавных персонажей, как Гринч, Лоракс и Хортон (см. [https://ru.wikipedia.org/wiki/Доктор\\_Сьюз](https://ru.wikipedia.org/wiki/Доктор_Сьюз)). — *Прим. пер.*

<sup>2</sup> Компьютерная сингулярность — точка во времени, с которой машины начинают совершенствоваться сами себя, без помощи кого-либо. — *Прим. пер.*

```
# жесткая пороговая функция в качестве активационной функции
def step_function(x):
    return 1 if x >= 0 else 0
```

```
# результат на выходе перцептрона
def perceptron_output(weights, bias, x):
    """возвращает 1, если перцептрон 'активизируется', и 0, если нет"""
    calculation = dot(weights, x) + bias
    return step_function(calculation)
```

Перцептрон попросту различает полупространства, разделенные гиперплоскостью точек  $x$ , для которых:

```
# скалярное произведение вектора весов и точек плюс смещение
dot(weights,x) + bias == 0
```

При правильно подобранных весах перцептрона могут решать ряд простых задач (рис. 18.1). Например, можно создать *логическую схему И* (которая возвращает 1, если оба входящих сигнала равны 1, и 0, если один из входящих сигналов равен 0):

```
weights = [2, 2] # вектор весов
bias = -3        # смещение
```

Если оба входящих сигнала равны 1, то взвешивание  $\text{calculation} = 2 + 2 - 3$  равно 1, и функция возвращает 1. Если только один из входящих сигналов равен 1, то  $\text{calculation} = 2 + 0 - 3 = -1$ , и функция возвращает 0. И если оба входящих сигнала равны 0, то  $\text{calculation} = -3$ , и функция возвращает 0.

Аналогичным образом создается *логическая схема ИЛИ*:

```
weights = [2, 2]
bias = -1
```

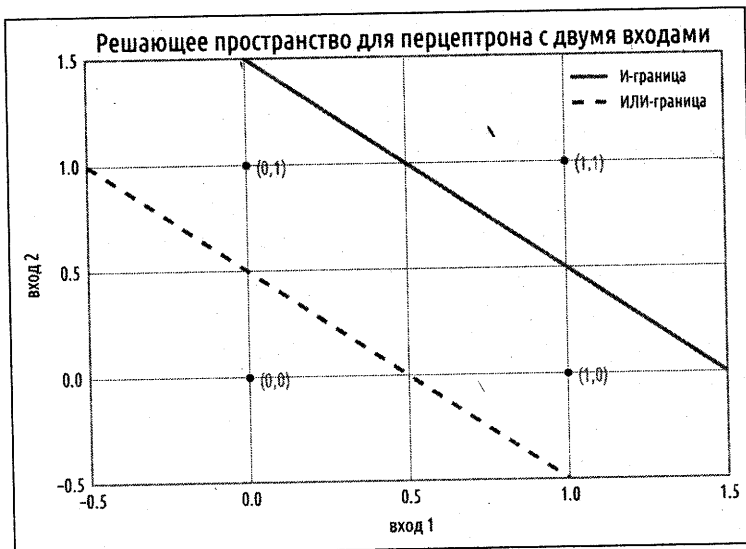


Рис. 18.1. Решающее пространство для перцептрона с двумя входами

Можно также создать *логическую схему НЕ* (которая имеет всего один входящий сигнал и преобразует 1 в 0 и наоборот):

```
weights = [-2]
bias = 1
```

Однако некоторые задачи однослойный перцептрон решить просто не способен. Например, как ни пытаться, но при помощи перцептрона не получится создать *логическую схему "исключающее ИЛИ"* (строгая конъюнкция), которая возвращает 1, если только один из его входящих сигналов равен 1, и 0 — в противном случае. Отсюда возникает потребность в более сложных моделях нейросети.

Естественно, для того чтобы создать логическую функцию, вовсе не нужно строить приближенную математическую модель нейрона:

```
and_gate = min      # булево И
or_gate = max       # булево ИЛИ
xor_gate = lambda x, y: 0 if x == y else 1 # булево исключающее ИЛИ
```

Как и настоящие нейроны, их искусственные аналоги вызывают интерес, только когда они соединены между собой.

## Нейронные сети прямого распространения

Человеческий мозг имеет чрезвычайно сложную топологию, и поэтому принято строить его приближенную модель на основе идеализированной нейронной сети *прямого распространения*, состоящей из дискретных *слоев* нейронов, где каждый слой присоединен к следующему. Такое строение, как правило, требует наличия слоя входов в сеть (который получает входящие сигналы и передает их дальше без изменений), один или несколько "скрытых слоев" (каждый из них состоит из нейронов, которые принимают выходящие сигналы предыдущего слоя, выполняют некие вычисления и передают результат в следующей слой) и слоя выходов из сети (который генерирует окончательные результаты).

Так же как в случае с перцептроном, каждый нейрон (не входного слоя) состоит из весов, которые поставлены в соответствие входящим сигналам, и смещения. Чтобы упростить представление сети, добавим смещение в конец вектора весов и установим для каждого нейрона *величину смещения* на входе, всегда равную 1.

Как и у перцептрона, будем суммировать произведение его входящих сигналов и весов. Однако здесь вместо пороговой активационной функции `step_function`, применяемой к произведению, будем возвращать ее гладкую аппроксимацию. В частности, воспользуемся сигмоидальной активационной функцией `sigmoid` (рис. 18.2):

```
# сигмоида в качестве функции активации
def sigmoid(t):
    return 1 / (1 + math.exp(-t))
```

Зачем применять сигмоидальную функцию `sigmoid` вместо более простой пороговой `step_function`? Для обучения нейронной сети используются дифференциальные



уравнения в частных производных, а для них требуются *гладкие* дифференцируемые функции. Пороговая функция даже не является непрерывной, а сигмоида представляет собой ее хорошую гладкую аппроксимацию.



Вы, наверное, обратили внимание на сигмоиду `sigmoid` в *главе 16*, где она носит название логистической `logistic`. Технически термин "сигмоида" относится к *форме* функции, а термин "логистический" — к данной конкретной функции, хотя часто эти термины используют как взаимозаменяемые.

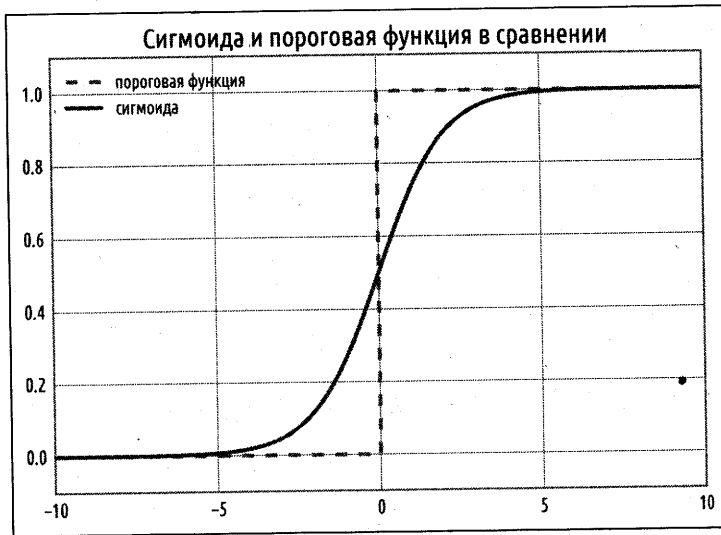


Рис. 18.2. Сигмоидальная и пороговая функции в сравнении

Затем вычислим выходящий сигнал нейрона:

```
# выходящий сигнал нейрона
def neuron_output(weights, inputs):
    return sigmoid(dot(weights, inputs))
```

При наличии этой функции нейрон может быть представлен простым списком весов, чья длина на единицу больше числа входов в этот нейрон (из-за веса смещения). Тогда нейронную сеть можно представить списком *слоев* (не входного уровня), где каждый слой — это просто список нейронов в этом слое.

Иными словами, нейронная сеть — это список (слоев) списков (нейронов) списков (весов).

В таком виде использовать нейронную сеть достаточно просто:

```
# механизм прямого распространения
def feed_forward(neural_network, input_vector):
    """принимает нейронную сеть (как список списков списков весов) и
    вектор входящих сигналов;
    возвращает результат прямого распространения входящих сигналов"""

    outputs = [] # выходы из сети
```

```

# обрабатывать послойно
for layer in neural_network:
    input_with_bias = input_vector + [1] # добавить величину смещения
    output = [neuron_output(neuron, input_with_bias) # вычислить результат
              for neuron in layer] # для этого нейрона
    outputs.append(output) # и запомнить его

# входом для следующего слоя становится
# вектор результатов текущего слоя
input_vector = output

return outputs

```

Теперь можно легко создать логическую схему для исключающего ИЛИ, которую не удалось реализовать на основе однослойного перцептрона. Для этого потребуются всего лишь скорректировать веса так, чтобы выходящие сигналы нейрона находились очень близко к 0 либо к 1:

```

xor_network = [# скрытый слой
               [[20, 20, -30], # нейрон 'И'
                [20, 20, -10]], # нейрон 'ИЛИ'
               # выходной слой
               [[-60, 60, -30]]] # нейрон '2-й, но не 1-й'

for x in [0, 1]:
    for y in [0, 1]:
        # функция feed_forward возвращает выходящие сигналы каждого нейрона
        # feed_forward[-1] - это выходящие сигналы нейронов слоя
        # выходов из сети
        print(x, y, feed_forward(xor_network, [x, y])[-1])

# 0 0 [9.38314668300676e-14]
# 0 1 [0.9999999999999059]
# 1 0 [0.9999999999999059]
# 1 1 [9.383146683006828e-14]

```

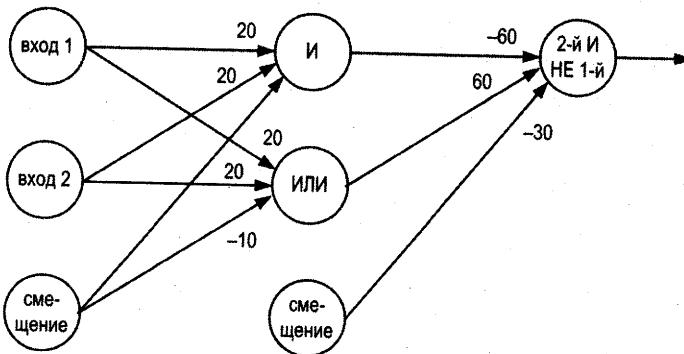


Рис. 18.3. Нейронная сеть логической схемы "исключающее ИЛИ"

При помощи скрытого слоя можно передать выходящий сигнал нейрона "И" и выходящий сигнал нейрона "ИЛИ" на вход нейрона "2-й, но не 1-й". В итоге получится сеть, которая выполняет операцию строгой дизъюнкции ("или, но не и"), соответствующую логической схеме "исключающее ИЛИ" (рис. 18.3).

## Метод обратного распространения ошибки

Как правило, нейронные сети не создают вручную, потому что с их помощью решают задачи гораздо более масштабные — к примеру, задача распознавания изображений может быть связана с применением сотен или даже тысяч нейронов, а также потому, что обычно невозможно "продумать до конца", какие значения должны генерировать нейроны.

Вместо этого нейронные сети *обучают* (как всегда) при помощи данных. Одним из популярных решений этой задачи является применение алгоритма, именуемого *методом обратного распространения ошибки*, который имеет определенную аналогию с алгоритмом градиентного спуска, рассмотренным ранее.

Допустим, дана обучающая выборка, которая состоит из входящих векторов и соответствующих целевых векторов, задающих требуемые результаты. Например, в предыдущем примере нейросети `xor_network` входящему вектору `[1, 0]` соответствовал целевой вектор `[1]`. Допустим также, что нейросеть располагает неким набором весов. Эти веса в дальнейшем корректируются на основе следующего алгоритма:

1. Выполнить прямой проход по сети, применив функцию `feed_forward` к входящему вектору, чтобы вычислить выходящие сигналы всех нейронов сети.
2. В результате в каждом выходном нейроне генерируется ошибка — разница между его выходящим и целевым значениями.
3. Вычислить градиент этой ошибки как функцию весов нейрона и затем скорректировать его веса в направлении, которое минимизирует ошибку.
4. Распространить ошибки выходящих сигналов на нейроны скрытого слоя, выполнив обратный проход по сети с вычислением ошибок для нейронов скрытого слоя.
5. Вычислить градиенты этих ошибок и скорректировать веса нейронов скрытого слоя так же, как в пункте 3.

Как правило, алгоритм выполняется большое число итераций для всей обучающей выборки до схождения сети (т. е. до тех пор, пока ошибка на всем множестве не достигнет приемлемого уровня):

```
# алгоритм обратного распространения ошибки
# на входе: сеть, вектор входов, вектор целей
def backpropagate(network, input_vector, targets):

    # выходы из скрытых слоев и выходы из сети
    hidden_outputs, outputs = feed_forward(network, input_vector)
```

```

# из производной сигмоидальной функции взято output * (1 - output)
output_deltas = [output * (1 - output) * (output - target)
                  for output, target in zip(outputs, targets)]

# понейронно скорректировать веса для слоя выходов (network[-1])
for i, output_neuron in enumerate(network[-1]):
    # взять i-й нейрон слоя выходов
    for j, hidden_output in enumerate(hidden_outputs + [1]):
        # скорректировать j-й вес на основе
        # дельты i-го нейрона и его j-го входящего значения
        output_neuron[j] -= output_deltas[i] * hidden_output

# распространить ошибки на скрытый слой, двигаясь в обратную сторону
hidden_deltas = [hidden_output * (1 - hidden_output) *
                  dot(output_deltas, [n[i] for n in output_layer])
                  for i, hidden_output in enumerate(hidden_outputs)]

# понейронно скорректировать веса для скрытого слоя (network[0])
for i, hidden_neuron in enumerate(network[0]):
    for j, input in enumerate(input_vector + [1]):
        hidden_neuron[j] -= hidden_deltas[i] * input

```

Это практически то же самое, если задать квадрат ошибки в качестве функции весов и использовать функцию `minimize_stochastic`, реализованную в *главе 8*.

В данном случае, расписать градиент в явном виде несколько проблематично. Для тех, кто знаком с математическим анализом и цепным правилом дифференцирования сложной функции, выкладки будут относительно простыми. Единственно озадачивает соблюдение правильности математических обозначений ("частная производная функции ошибки относительно веса, который нейрон  $i$  назначает входящему значению, приходящему из нейрона  $j$ ).

## Пример: преодоление капчи

Чтобы лица, которые регистрируются на сайте социальной сети, гарантированно были людьми, а не компьютерными ботами, директор по управлению продуктами хочет внедрить в процесс регистрации посетителей *капчу*<sup>3</sup>. В частности, он хотел бы показывать пользователю изображения цифр и предлагать ему вводить цифры для подтверждения, что он или она — человек.

Он не верит, что компьютеры способны легко справиться с этой задачей, и поэтому вы собираетесь переубедить его, создав программу, которая эту задачу решает.

<sup>3</sup> Капча (от CAPTCHA, Completely Automated Public Turing test to tell Computers and Humans Apart — полностью автоматизированный публичный тест Тьюринга для различения компьютеров и людей) — компьютерный тест, используемый для того, чтобы определить, кем является пользователь системы: человеком или компьютером (см. <https://ru.wikipedia.org/wiki/Капча>). — Прим. пер.

Представим каждую цифру в виде изображения 5×5 пикселей:

```

e...e  ..e..  e...e  e...e  e...e  e...e  e...e  e...e  e...e  e...e
e...e  ..e..  ....e  ....e  e...e  e...e  e...e  ....e  e...e  e...e
e...e  ..e..  e...e  e...e  e...e  e...e  e...e  ....e  e...e  e...e
e...e  ..e..  e...e  ....e  ....e  ....e  e...e  ....e  e...e  ....e
e...e  ..e..  e...e  e...e  ....e  e...e  e...e  ....e  e...e  e...e

```

Нейронная сеть требует, чтобы на вход поступал вектор чисел, поэтому развернем каждое изображение в вектор длиной 25 элементов, где каждый элемент равен либо 1 ("пиксел установлен"), либо 0 ("пиксел обнулен").

Например, цифра 0 будет представлена так:

```

zero_digit = [1,1,1,1,1,
              1,0,0,0,1,
              1,0,0,0,1,
              1,0,0,0,1,
              1,0,0,0,1,
              1,1,1,1,1]

```

На выходе нейросеть должна указывать на предполагаемую цифру. Для этого понадобится 10 выходных значений. К примеру, правильным результатом для цифры 4 будет:

```
[0, 0, 0, 0, 1, 0, 0, 0, 0, 0]
```

Тогда, при условии, что входы `inputs` правильно упорядочены от 0 до 9, целыми будут:

```

targets = [[1 if i == j else 0 for i in range(10)] # диагон. матрица 10x10
           for j in range(10)]

```

благодаря чему (к примеру) значение `targets[4]` соответствует правильному выходному значению для цифры 4. На данный момент все готово, чтобы построить нейронную сеть:

```

random.seed(0) # это необходимо для получения повторяющихся результатов
input_size = 25 # каждое входящее значение - это вектор
               # длиной 25 элементов
num_hidden = 5 # 5 нейронов в скрытом слое
output_size = 10 # 10 значений на выходе для каждого входящего

# каждый скрытый нейрон имеет один вес для входного нейрона
# плюс вес смещения
hidden_layer = [[random.random() for __ in range(input_size + 1)]
                for __ in range(num_hidden)]

# каждый выходной нейрон имеет один вес для скрытого нейрона
# плюс вес смещения
output_layer = [[random.random() for __ in range(num_hidden + 1)]
                for __ in range(output_size)]

```

```
# сеть начинает работу со случайными весами
network = [hidden_layer, output_layer]
```

Обучим ее при помощи алгоритма обратного распространения ошибки:

```
# для схождения сети 10 000 итераций, по-видимому, будет достаточно
for __ in range(10000):
    for input_vector, target_vector in zip(inputs, targets):
        backpropagate(network, input_vector, target_vector)
```

Сеть хорошо справляется с обучающими примерами, что очевидно из результатов:

```
# функция прогнозирования
def predict(input):
    return feed_forward(network, input)[-1]
```

```
predict(inputs[7])
# [0.026, 0.0, 0.0, 0.018, 0.001, 0.0, 0.0, .967, 0.0, 0.0]
```

Результат показывает, что нейрон выхода для цифры 7 генерирует 0.97, тогда как другие нейроны показывают очень малые значения.

Помимо этого, нейросеть справляется с цифрами, нарисованными по-другому, такими как эта стилизованная цифра 3:

```
predict([0,1,1,1,0, # .@@@.
        0,0,0,1,1, # ...@@
        0,0,1,1,0, # ..@@.
        0,0,0,1,1, # ...@@
        0,1,1,1,0]) # :@@@.

# [0.0, 0.0, 0.0, 0.92, 0.0, 0.0, 0.0, 0.01, 0.0, 0.12]
```

Нейросеть продолжает считать, что эта цифра похожа на тройку, а вот в отношении стилизованной цифры 8 мнение разделилось между 5, 8 и 9:

```
predict([0,1,1,1,0, # .@@@.
        1,0,0,1,1, # @..@@
        0,1,1,1,0, # .@@@.
        1,0,0,1,1, # @..@@
        0,1,1,1,0]) # .@@@.

# [0.0, 0.0, 0.0, 0.0, 0.0, 0.55, 0.0, 0.0, 0.93, 1.0]
```

При наличии более объемной обучающей выборки с подобной ситуацией, наверное, удалось бы справиться.

Несмотря на то, что нейросеть выполняет работу не совсем прозрачно, тем не менее, есть возможность исследовать веса скрытого слоя и в результате получить представление о том, что в действительности они распознают. В частности, можно нарисовать веса каждого нейрона в виде матрицы размера  $5 \times 5$ , соответствующей входам  $5 \times 5$ .

В реальной ситуации, вероятно, было бы лучше изобразить нулевой вес белым цветом, более крупные положительные веса различными оттенками (к примеру) зеленого, а более крупные отрицательные веса оттенками (допустим) красного. К сожалению, это трудно сделать в формате черно-белой книги.

Вместо этого изобразим нулевые веса белым, веса, более удаленные от него, сделаем более темными, а для обозначения отрицательных весов применим штриховку.

Для этих целей воспользуемся функцией `pyplot.imshow`, которая ранее еще не встречалась. С ее помощью можно создавать изображения, рисуя пиксел за пикселом. В обычных условиях при решении аналитических задач она не особо практически полезна, однако тут нужен именно такой функционал<sup>4</sup>:

```
import matplotlib
# веса
weights = network[0][0]          # первый нейрон в скрытом слое
abs_weights = map(abs, weights)  # темный цвет зависит
                                # от абсолютного значения

grid = [abs_weights[row:(row+5)] # преобразовать веса в матрицу 5 x 5
         for row in range(0,25,5)] # [weights[0:5], ..., weights[20:25]]

ax = plt.gca()                  # для штриховки нужен объект axis (ось)

ax.imshow(grid,                 # здесь так же, как в plt.imshow
           cmap=matplotlib.cm.binary, # использовать шкалу ч/б цвета
           interpolation='none')    # рисовать блоками

def patch(x, y, hatch, color):
    """вернуть объект 'patch' библиотеки matplotlib с указанными
    координатами, шаблоном штриховки и цветом"""
    return matplotlib.patches.Rectangle((x - 0.5, y - 0.5), 1, 1,
                                         hatch=hatch, fill=False, color=color)

# заштриховать отрицательные веса
for i in range(5):              # строка
    for j in range(5):          # столбец
        if weights[5*i + j] < 0: # строка i, столбец j = weights[5*i + j]
            # добавить ч/б штриховку, чтобы было видно на темном или светлом
            ax.add_patch(patch(j, i, '/', "white"))
            ax.add_patch(patch(j, i, '\\', "black"))

plt.show()
```

<sup>4</sup> При работе с ниже приведенным кодом в Python 3 для получения вектора абсолютных весов следует скорректировать выражение так: `abs_weights = list(map(abs, weights))`. Кроме того, окончательные значения смещений могут незначительно отличаться от значений, приведенных в книге. — Прим. пер.

На рис. 18.4 видно, что первый скрытый нейрон имеет большие положительные веса в левом столбце и в центре средней строки, и большие отрицательные веса в правом столбце. (К тому же можно убедиться, что у него довольно большое отрицательное смещение, а значит, он сработает только, если или когда на входе получит именно те положительные значения, которые "ожидает").

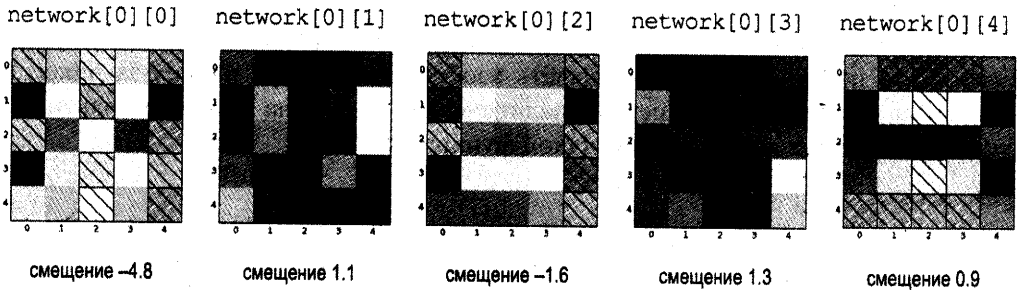


Рис. 18.4. Веса для скрытого слоя

Действительно, при таких входящих значениях сеть делает то, что от нее ожидается:

```
# только левый столбец
```

```
left_column_only = [1, 0, 0, 0, 0] * 5
```

```
print(feed_forward(network, left_column_only)[0][0]) # 1.0
```

```
# центр, средняя строка
```

```
center_middle_row = [0, 0, 0, 0, 0] * 2 + [0, 1, 1, 1, 0] + [0, 0, 0, 0, 0] * 2
```

```
print(feed_forward(network, center_middle_row)[0][0]) # 0.95
```

```
# только правый столбец
```

```
right_column_only = [0, 0, 0, 0, 1] * 5
```

```
print(feed_forward(network, right_column_only)[0][0]) # 0.0
```

Аналогичным образом, среднему скрытому нейрону, похоже, "нравятся" горизонтальные линии, в отличие от боковых вертикальных линий. А последнему скрытому нейрону — центральная строка, в отличие от правого столбца. (Два других нейрона интерпретировать труднее.)

Что произойдет, если пропустить через нейросеть стилизованную цифру 3?

```
my_three = [0,1,1,1,0, # .@@@.
            0,0,0,1,1, # ...@@
            0,0,1,1,0, # ..@@.
            0,0,0,1,1, # ...@@
            0,1,1,1,0] # .@@@.
```

```
hidden, output = feed_forward(network, my_three)
```

Выходные значения hidden будут следующими:

```
0.121080 # из network[0][0], вероятно, навечно (1, 4)
0.999979 # из network[0][1], большой вклад от (0, 2) и (2, 2)
```



```
0.999999 # из network[0][2], положительно везде, кроме (3, 4)
0.999992 # из network[0][3], снова большой вклад от (0, 2) и (2, 2)
0.000000 # из network[0][4], отрицательно или 0 везде,
           # кроме центральной строки
```

которые входят в нейрон выхода "три" с весами `network[-1][3]`:

```
-11.61 # вес для hidden[0]
-2.17 # вес для hidden[1]
 9.31 # вес для hidden[2]
-1.38 # вес для hidden[3]
-11.47 # вес для hidden[4]
-1.92 # вес для смещения на входе
```

благодаря чему нейрон вычислит:

```
sigmoid(.121 * -11.61 + 1 * -2.17 + 1 * 9.31 - 1.38 * 1 - 0 * 11.47 - 1.92)
```

Это в результате даст 0.92, как и было. В сущности, скрытый слой вычисляет пять разных сегментов 25-мерного пространства, отображая каждое входящее 25-мерное значение в пять чисел. И затем каждый нейрон выхода смотрит только на результаты этих пяти сегментов.

Как можно было убедиться, `my_three` слегка "склоняется" в сторону сегмента 0 (т. е. незначительно активирует скрытый нейрон 0), сильно "склоняется" в сторону сегментов 1, 2 и 3 (т. е. сильно активирует эти скрытые нейроны) и почти "равнодушен" к сегменту 4 (т. е. не активирует тот нейрон вообще).

Затем каждый из 10 нейронов выхода из сети использует только эти пять активаций, чтобы решить, является ли `my_three` их цифрой или нет.

## Для дальнейшего изучения

- ◆ Образовательная онлайн-платформа Coursera предлагает бесплатный курс по теме "Искусственные нейронные сети в машинном обучении" (<https://www.coursera.org/course/neuralnets>). В последний раз он проводился в 2012 г., однако материалы курса по-прежнему доступны.
- ◆ Майкл Нильсен (Michael Nielsen) пишет бесплатную онлайн-книгу на тему "Нейронные сети и глубокое обучение" (<http://neuralnetworksanddeeplearning.com/>). К тому времени, когда вы будете читать это, он должен уже ее завершить.
- ◆ PyBrain (<http://pybrain.org/>) — достаточно простая нейросетевая библиотека, написанная на языке Python.
- ◆ Pylearn2 (<http://deeplearning.net/software/pylearn2/>) — намного более продвинутая (но и более сложная в использовании) нейросетевая библиотека.

# Кластеризация

Где славный наш союз  
Мятежно избавлял от уз.  
— Роберт Геррик<sup>1</sup>

Большинство алгоритмов в этой книге относятся к классу контролируемого обучения или обучения с учителем, т. е. они начинают работу с набора *маркированных* данных и используют их в качестве основы для выполнения предсказаний в отношении новых, *немаркированных* данных. Кластеризация (clustering) — это пример неконтролируемого обучения или обучения без учителя, когда работа ведется с совершенно *немаркированными* данными (или же если данные промаркированы, то метки игнорируются).

## Идея

Обратившись к некоему источнику данных, можно заметить, что данные, скорее всего, каким-то образом собраны в группы или *кластеры* (термины синонимичны). Набор данных, демонстрирующий, где живут миллионеры, вероятно, сформирует группы в таких местах, как Беверли Хиллз и Манхэттен. Набор данных, показывающий, сколько часов в неделю люди проводят на работе, вероятно, создаст группу вокруг 40 часов (а если данные взяты из штата с законами, которые предусматривают специальные льготы для людей, работающих не более 20 часов в неделю, то он, скорее всего, создаст еще одну группу непосредственно около 19). Набор данных о демографии зарегистрированных избирателей, скорее всего, сформирует несколько групп (например, "футбольные мамочки", "скучающие пенсионеры", "безработные представители поколения Y"), которые социологи и политические консультанты, по-видимому, посчитают актуальными.

В отличие от большинства задач, которые были рассмотрены ранее, "правильной" кластеризации не существует. Альтернативная схема кластеризации может сгруппировать некоторых "безработных поколения Y" с "аспирантами", другая — с "обитателями родительских подвалов"<sup>2</sup>. Ни одна схема не является неизбежно

<sup>1</sup> Роберт Геррик (1591–1674) — английский поэт, представитель группы так называемых "поэтов-кавалеров", сторонников короля Карла I. Здесь обыгрывается слово "союз" (cluster). — *Прим. пер.*

<sup>2</sup> Обитатели родительских подвалов (<https://resurgencemedia.net/2015/10/19/millennial-basement-dwellers-and-their-pathetic-parents/>) — термин, чаще всего относящийся к иждивенцам у родителей. Это люди, которые проводят всю или большую часть времени в подвале дома. Проживание в подвале — это форма ухода от действительности и современный тренд в США, связанный с использованием последних достижений в области электроники и вычислительной техники, но который средни тенденции

правильнее другой — напротив, каждая, скорее всего, является оптимальнее относительно своей метрики, измеряющей качество кластеризации.

Кроме того, группы сами себя не маркируют. Вам придется делать это самостоятельно, проверяя данные, расположенные внутри каждой из них.

## Модель

В данном конкретном случае каждое входящее значение `input` — это вектор в  $d$ -мерном пространстве (который, как обычно, будет представлен в виде списка чисел). Задача будет заключаться в том, чтобы определять группы однотипных входящих значений и (иногда) находить репрезентативное значение для каждой группы.

Например, каждое входящее значение может быть заголовком поста в сетевом журнале (в виде числового вектора, который каким-то образом этот пост представляет), и в этом случае задача может состоять в том, чтобы найти группы похожих постов, возможно, для того, чтобы понять, о чем пользователи сайта пишут в своих журналах. Или, предположим, некая фотография содержит тысячи цветов (в виде триплетов, состоящих из красного, зеленого, синего), и нужно сделать ее экранную копию в 10-цветном формате. Здесь кластеризация способна помочь выбрать 10 цветов и при этом свести к минимуму суммарную "ошибку уменьшения цвета"<sup>3</sup>.

Одним из простейших методов кластеризации или кластерного анализа является *метод  $k$  средних*, в котором число кластеров  $k$  выбирается заранее. Далее задача состоит в том, чтобы сегментировать входящие значения на подмножества  $S_1, \dots, S_k$  таким образом, чтобы минимизировать полную сумму квадратов расстояний от каждой точки до среднего значения назначенного ей кластера.

Назначить  $n$  точек данных  $k$  группам можно самыми разными способами, вследствие чего найти оптимальную конфигурацию групп очень трудно. Мы остановимся на итеративном алгоритме, который обычно находит хорошую конфигурацию.

1. Начать с множества  $k$  средних, т. е. точек в  $d$ -мерном пространстве.
2. Назначить каждую точку ближайшему среднему значению.
3. Если ни у одной точки ее назначение не изменилось, то остановиться и сохранить группы.
4. Если назначение одной из точек изменилось, то пересчитать средние значения и вернуться к шагу 2.

При помощи функции `vector_mean` из главы 4 (которая возвращает вектор средних значений списка векторов) достаточно просто создать питоновский класс, который это выполняет:

---

предыдущих поколений, когда такой выбор делался большей частью из желания избежать очного контакта с другими людьми. — Прим. пер.

<sup>3</sup> <http://www.imagemagick.org/script/quantize.php#measure>.

```

class KMeans:
    """класс выполняет кластеризацию по методу k средних"""

    def __init__(self, k):
        self.k = k          # число кластеров
        self.means = None  # средние кластеров

    def classify(self, input):
        """вернуть индекс кластера, ближайшего к входящему значению input"""
        return min(range(self.k),
                    key=lambda i: squared_distance(input, self.means[i]))

    def train(self, inputs):
        # выбрать k случайных точек в качестве исходных средних
        self.means = random.sample(inputs, self.k)
        assignments = None
        while True:
            # найти новые назначения
            new_assignments = map(self.classify, inputs)

            # если ни одно назначение не изменилось, то завершить
            if assignments == new_assignments:
                return

            # в противном случае сохранить новые назначения
            assignments = new_assignments

            # и вычислить новые средние на основе новых назначений
            for i in range(self.k):
                # найти все точки, назначенные кластеру i
                i_points = [p for p, a in zip(inputs, assignments) if a == i]

                # удостовериться, что i_points не пуст, чтобы не делить на 0
                if i_points:
                    self.means[i] = vector_mean(i_points)

```

Посмотрим, как это работает.

## Пример: встречи для специалистов

Чтобы отметить рост пользовательской базы соцсети DataSciencester, директор по премированию клиентуры хочет организовать для пользователей, проживающих в одном городе, несколько встреч ИТ-специалистов (meetups) с пивом, пиццей и фирменными футболками DataSciencester. Вы знаете места проживания всех местных пользователей (рис. 19.1), и директор хотел бы, чтобы вы сами определили места встреч, которые были бы удобными для всех желающих приехать.

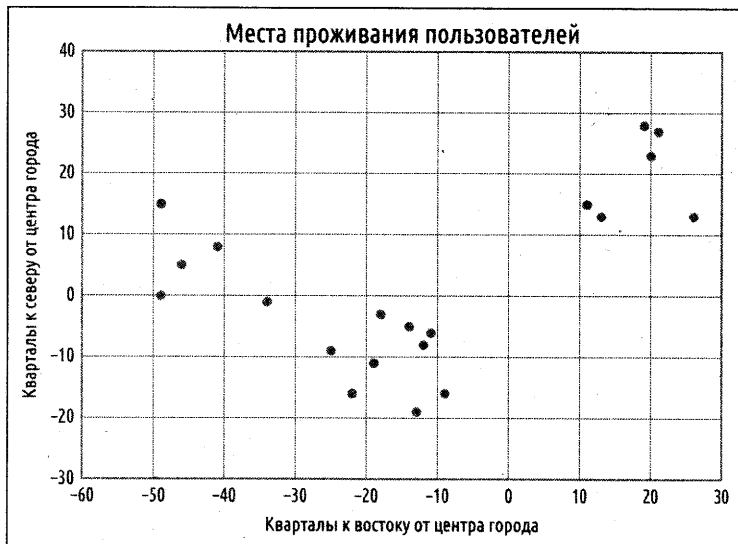


Рис. 19.1. Места проживания пользователей в городе

В зависимости от того, как смотреть, можно увидеть два или три кластера. (В этом легко убедиться визуально, поскольку данные представлены всего лишь в двух измерениях. В случае большего числа измерений было бы намного труднее это заметить.)

Для начала представим, что выделен бюджет, достаточный для трех встреч. Вы садитесь за компьютер и пробуете следующее<sup>4</sup>:

```
random.seed(0) # чтобы были повторимые результаты
clusterer = KMeans(3)
clusterer.train(inputs)
print(clusterer.means)
```

Вы находите три группы с центрами в районах  $[-44, 5]$ ,  $[-16, -10]$  и  $[18, 20]$ , подбираете места встреч рядом с этими районами (рис. 19.2) и показываете их директору. Тот сообщает, что бюджета, к сожалению, хватит всего на две встречи. "Нет проблем", — говорите вы:

```
random.seed(0)
clusterer = KMeans(2)
clusterer.train(inputs)
print(clusterer.means)
```

Как показано на рис. 19.3, одну встречу следует провести рядом с районом  $[18, 20]$ , а другую теперь рядом с районом  $[-26, -5]$ .

<sup>4</sup> В Python 3.5.2 результаты, повторимые с приведенными в книге, достигаются при удалении оператора `random.seed(0)`. — Прим. пер.

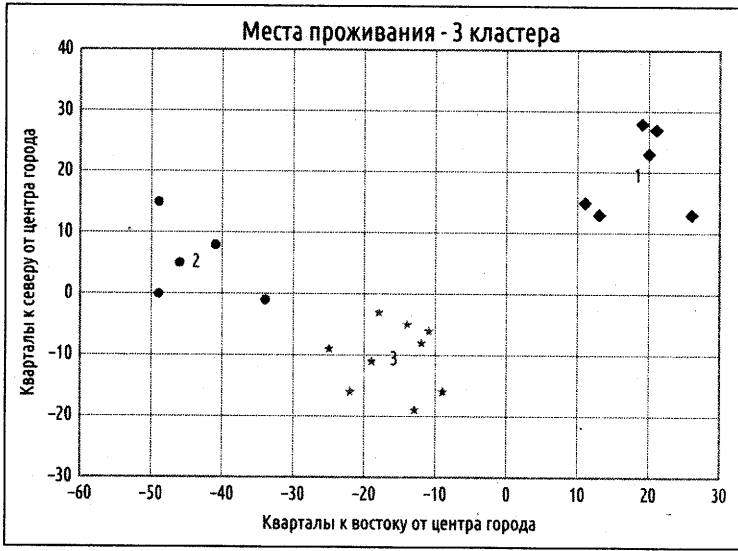


Рис. 19.2. Места проживания пользователей, сгруппированные в три кластера

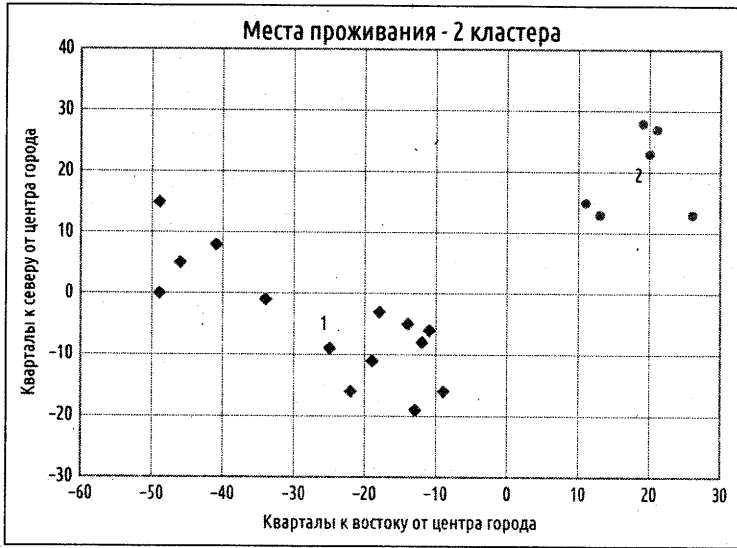


Рис. 19.3. Места проживания пользователей, сгруппированные в два кластера

## Выбор числа $k$

В предыдущем примере, выбор числа  $k$  был обусловлен факторами, находящимися вне вашего контроля. Обычно так не бывает. Существует целый ряд способов выбора числа  $k$ . Один из самых легких для понимания предполагает построение диаграммы суммарного квадратичного отклонения (между каждой точкой и средним своей группы), как функции от  $k$ , и нахождение места "излома" кривой:

```

# суммарное квадратичное отклонение кластеризации
def squared_clustering_errors(inputs, k):
    """находит суммарное квадратичное отклонение (ошибку)
    от k средних при кластеризации входящих данных"""
    clusterer = KMeans(k)
    clusterer.train(inputs)
    means = clusterer.means
    assignments = map(clusterer.classify, inputs)

    return sum(squared_distance(input, means[cluster])
               for input, cluster in zip(inputs, assignments))

# теперь нарисовать от 1 до len(inputs) кластеров

ks = range(1, len(inputs) + 1)
errors = [squared_clustering_errors(inputs, k) for k in ks]

plt.plot(ks, errors)
plt.xticks(ks)
plt.xlabel("k")
plt.ylabel("Суммарное квадратичное отклонение")
plt.title("Суммарная ошибка и число кластеров")
plt.show()

```

График на рис. 19.4 наглядно демонстрирует, что этот метод согласуется с первоначальной визуальной оценкой о "правильности" числа групп, равного 3.

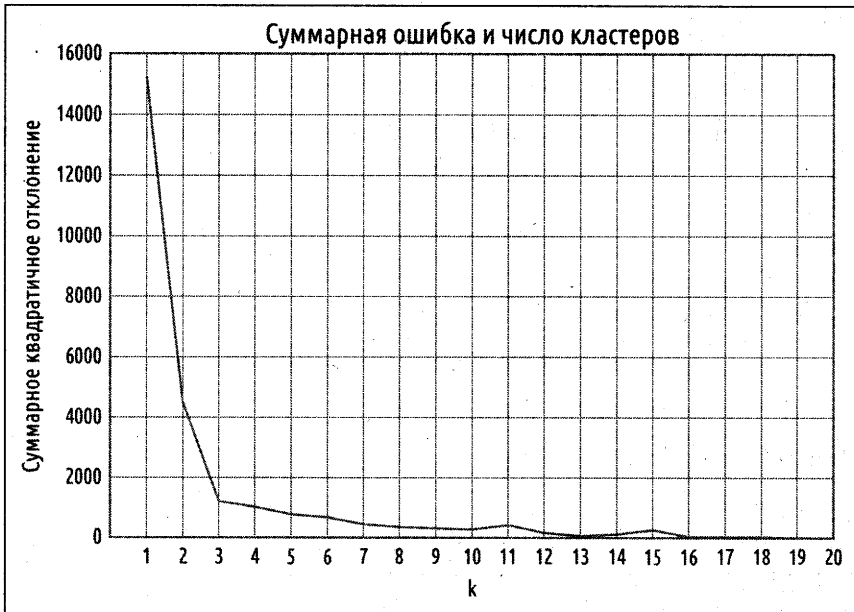


Рис. 19.4. Выбор k

## Пример: кластеризация цвета

Директор по промоутерским материалам разработал привлекательные фирменные наклейки DataSciencester, которые он хотел бы раздать на дружеских встречах. К сожалению, офисный стикер-принтер может печатать только в пяти цветах, а поскольку директор по искусству находится в творческом отпуске, он интересуется, нельзя ли как-нибудь изменить дизайн наклеек так, чтобы они содержали только пять цветов.

Компьютерные изображения можно представить в виде двумерного массива пикселей, где каждый пиксел сам является трехмерным вектором (красный, зеленый, синий), обозначающим его цвет.

Тогда создание пятицветной версии изображения подразумевает:

- ◆ выбор пяти цветов;
- ◆ назначение одного из этих цветов каждому пикселу.

Оказывается, что эта задача легко решается на основе кластеризации методом  $k$  средних, который способен сегментировать пиксели на 5 групп в красно-зелено-синем пространстве. Если затем переокрасить пиксели в каждой группе в средний цвет, то задача будет решена.

Для начала надо загрузить изображение в Python. Как оказалось, это можно сделать из библиотеки `matplotlib`:

```
path_to_png_file = r"C:\images\image.png" # указать путь к изображению
import matplotlib.image as mpimg
img = mpimg.imread(path_to_png_file)      # считать изображение
```

В основе изображения `img` лежит массив библиотеки NumPy, но в целях изложения будем рассматривать его как список списков списков.

`img[i][j]` — это пиксел в  $i$ -й строке и  $j$ -м столбце, причем каждый пиксел — это список `[red, green, blue]` чисел в диапазоне от 0 до 1, которые обозначают цвет этого пиксела<sup>5</sup>:

```
top_row = img[0]
top_left_pixel = top_row[0]
red, green, blue = top_left_pixel
```

В частности, развернутый список всех пикселей можно получить так:

```
pixels = [pixel for row in img for pixel in row]
```

и затем передать его в кластеризатор:

```
clusterer = KMeans(5)
clusterer.train(pixels) # это может потребовать некоторого времени
```

<sup>5</sup> [https://en.wikipedia.org/wiki/RGB\\_color\\_model](https://en.wikipedia.org/wiki/RGB_color_model).



Когда кластеризатор завершит работу, создадим новое изображение в том же формате:

```
def recolor(pixel):
    cluster = clusterer.classify(pixel) # индекс ближайшего кластера
    return clusterer.means[cluster]    # среднее ближайшего кластера

new_img = [[recolor(pixel) for pixel in row] # переокрасить эту строку
           # пикселей
           for row in img] # для каждой строки в изображении
```

и выведем его на экран при помощи `plt.imshow()`:

```
plt.imshow(new_img)
plt.axis('off')
plt.show()
```

Трудно изобразить цветные результаты в черно-белой книге, впрочем, фотографии на рис. 19.5 показывают полутоновую версию полноцветной фотографии и результат после понижения шкалы до пяти цветов.



Рис. 19.5. Оригинальная фотография и обесцвеченная по методу 5 средних

## Восходящий метод иерархической кластеризации

Альтернативный подход к кластеризации заключается в "выращивании" кластеров снизу вверх. Это делается следующим образом:

1. Сделать каждое входящее значение одноточечным кластером (из одного элемента).
2. Пока остается больше одного кластера, найти два ближайших и объединить.

В итоге получится один гигантский кластер, содержащий все входящие значения. Если отслеживать порядок объединения, то можно воссоздать любое число кластеров путем их разъединения. Например, если требуется три кластера, то надо всего лишь отменить последние два объединения.

Применим для кластеров очень простой вид представления. Значения будут размещаться в *листовых* кластерах в виде одноэлементных кортежей.

```
# чтобы получить одноэлементный кортеж, нужно поставить в конце запятую;
# иначе Python примет скобку за скобку
leaf1 = ([10, 20],)
leaf2 = ([30, -15],)
```

Они используются для создания *объединений* кластеров, которые в свою очередь представлены двухэлементными кортежами в формате (порядок объединения, дочерние элементы):

```
merged = (1, [leaf1, leaf2])
```

Обсудим порядок объединения чуть позже, а пока создадим несколько вспомогательных функций:

```
# лиственной кластер?
def is_leaf(cluster):
    """кластер является листом, если его длина = 1"""
    return len(cluster) == 1

# получить дочерние элементы
def get_children(cluster):
    """вернуть два дочерних элемента данного кластера,
    если он объединенный кластер; вызывает исключение,
    если это лиственной кластер"""
    if is_leaf(cluster):
        raise TypeError("лиственной кластер не имеет дочерних элементов")
    else:
        return cluster[1]

# получить значения
def get_values(cluster):
    """вернуть значение в кластере (если это лиственной кластер)
    или все значения в лиственных кластерах под ним (если нет)"""
    if is_leaf(cluster):
        return cluster # это уже одноэлементный кортеж,
                        # содержащий значение
    else:
        return [value
                for child in get_children(cluster)
                for value in get_values(child)]
```

Чтобы выполнить объединение ближайших кластеров, требуется некоторое представление о расстоянии между кластерами. Воспользуемся *минимальным расстоянием*

нием между элементами двух кластеров, при котором объединяются два кластера, расположенных максимально близко (что иногда приводит к большим кластерам с не очень плотной цепеобразной структурой). Если нужно получить плотные сферические кластеры, то вместо этого можно воспользоваться *максимальным* расстоянием, поскольку оно объединяет два кластера, которые помещаются в самый маленький клубок. Оба варианта одинаково распространены наравне со *средним* расстоянием (центроидный метод):

```
# расстояние между двумя кластерами
def cluster_distance(cluster1, cluster2, distance_agg=min):
    """вычислить все попарные расстояния между cluster1 и cluster2
    и применить функцию объединения _distance_agg
    к результирующему списку"""
    return distance_agg([distance(input1, input2)
                          for input1 in get_values(cluster1)
                          for input2 in get_values(cluster2)])
```

Элемент кортежа с порядковым номером объединения необходим для отслеживания порядка, в котором выполняются объединения кластеров. Меньшее число будет обозначать *более позднее* объединение. Вследствие этого, разъединение кластеров выполняется от наименьшего значения порядка объединения к наибольшему. Поскольку листовые кластеры не объединяются (и, значит, не подлежат разъединению), то присвоим им литерал положительной бесконечности:

```
# получить порядковый номер объединения
def get_merge_order(cluster):
    if is_leaf(cluster):
        return float('inf')
    else:
        # порядковый номер merge_order - это
        return cluster[0] # 1-й элемент 2-элементного кортежа
```

Теперь все готово для создания алгоритма кластеризации:

```
# восходящий (снизу вверх) алгоритм кластеризации
def bottom_up_cluster(inputs, distance_agg=min):
    # начать с того, что все входы - листовые кластеры / 1-элементный кортеж
    clusters = [(input,) for input in inputs]

    # пока остается более одного кластера...
    while len(clusters) > 1:
        # найти два ближайших кластера
        c1, c2 = min([(cluster1, cluster2)
                      for i, cluster1 in enumerate(clusters)
                      for cluster2 in clusters[:i]],
                     key=lambda (x, y): cluster_distance(x, y, distance_agg))

        # исключить их из списка кластеров
        clusters = [c for c in clusters if c != c1 and c != c2]
```



18 на основе сочетания листьев [19, 28] и [21, 27]. В результате последнего объединения получен кластер 0.

Тем не менее, обычно стараются избегать подобного рода заставляющих шуриться текстовых представлений. (Впрочем, создание удобной в использовании визуализации кластерной иерархии могло бы стать увлекательным упражнением.) Вместо этого напишем функцию, которая генерирует любое количество кластеров, выполнив положенное количество разъединений:

```
# сгенерировать кластеры
def generate_clusters(base_cluster, num_clusters):
    # начать со списка, состоящего только из базового кластера
    clusters = [base_cluster]

    # продолжать, пока кластеров не достаточно...
    while len(clusters) < num_clusters:
        # выбрать из кластеров тот, который был объединен последним
        next_cluster = min(clusters, key=get_merge_order)
        # исключить его из списка
        clusters = [c for c in clusters if c != next_cluster]
        # и добавить его дочерние элементы к списку
        # (т. е. разъединить его)
        clusters.extend(get_children(next_cluster))

    # когда уже достаточно кластеров...
    return clusters
```

Например, если нужно сгенерировать три кластера, то можно поступить так:

```
three_clusters = [get_values(cluster)
                  for cluster in generate_clusters(base_cluster, 3)]
```

и затем вывести их на диаграмму:

```
for i, cluster, marker, color in zip([1, 2, 3],
                                     three_clusters,
                                     ['D', 'o', '*'],
                                     ['r', 'g', 'b']):
    xs, ys = zip(*cluster) # трюк с разъединением списка
    plt.scatter(xs, ys, color=color, marker=marker)

    # установить число в среднее значение кластера
    x, y = vector_mean(cluster)
    plt.plot(x, y, marker='$' + str(i) + '$', color='black')

plt.title("Места проживания - 3 кластера снизу вверх, min")
plt.xlabel("Кварталы к востоку от центра города")
plt.ylabel("Кварталы к северу от центра города")
plt.show()
```



Рис. 19.6. Три кластера по восходящему методу кластеризации на основе функции min

Это даст результаты, сильно отличающиеся от тех, которые были получены на основе метода  $k$  средних, как показано на рис. 19.6.

Как уже упоминалось ранее, это происходит потому, что применение целевой функции min для вычисления расстояния между кластерами в `cluster_distance` имеет свойство порождать цепеобразные кластеры. Если вместо нее использовать целевую функцию max (которая порождает плотные кластеры), то результат будет выглядеть точно так же, как и в случае с 3 средними (рис. 19.7).

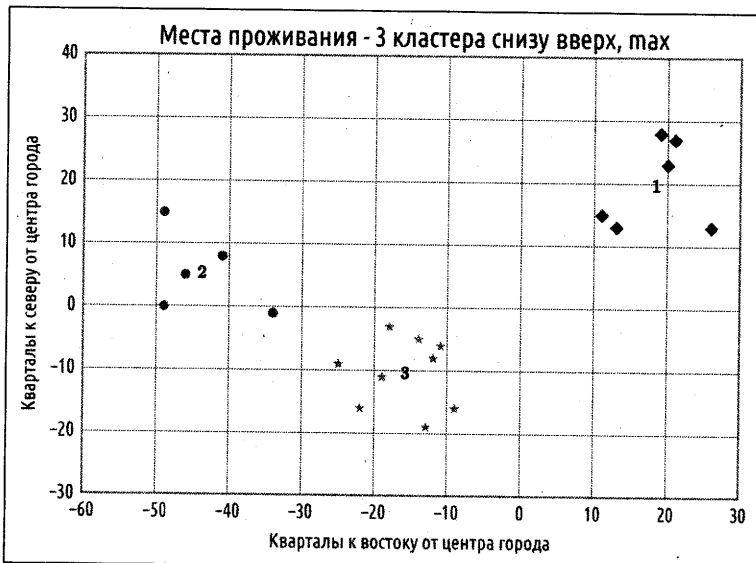


Рис. 19.7. Три кластера по восходящему методу кластеризации на основе функции max



Показанная выше реализация восходящего метода кластеризации `bottom_up_clustering` относительно проста и одновременно крайне неэффективна. В частности, здесь на каждом шаге заново вычисляется расстояние между каждой парой входящих значений. В более эффективной реализации попарные расстояния будут вычисляться заранее и затем будет выполняться поиск внутри списка расстояний `cluster_distance`. В действительно эффективной реализации, скорее всего, также будут учитываться расстояния `cluster_distance` предыдущего шага.

## Для дальнейшего изучения

- ◆ Библиотека `scikit-learn` содержит модуль `sklearn.cluster` (<http://scikit-learn.org/stable/modules/clustering.html>), который располагает несколькими алгоритмами кластеризации, включая алгоритм  $k$  средних `KMeans` и алгоритм иерархической кластеризации Ворда `Ward` (в котором используется критерий объединения кластеров, отличающийся от рассмотренного в этой главе).
- ◆ Библиотека `SciPy` (<http://www.scipy.org/>) располагает двумя моделями кластеризации `scipy.cluster.vq` (которая работает по методу  $k$  средних) и `scipy.cluster.hierarchy` (которая содержит несколько разновидностей алгоритма иерархической кластеризации).

# Обработка естественного языка

Они оба побывали на каком-то словесном пиру и наворовали там объедков.

*Уильям Шекспир<sup>1</sup>*

*Обработка естественного языка* (natural language processing) имеет отношение к вычислительным технологиям с участием языка. Они представляют собой обширнейшую область, однако мы обратимся лишь к нескольким приемам, простым и не очень.

## Облака слов

В главе 1 были вычислены частотности слов, встречающихся в темах, которые интересуют пользователей. Один из приемов визуализации слов и их частотностей представлен облаками слов (облаками тегов или визуализациями частотности слов в виде взвешенного списка), которые эффектно расставляют слова в соответствии с размерами, пропорциональными их частотностям.

Надо отметить, что аналитики данных обычно не высокого мнения об облаках слов во многом потому, что размещение слов не означает ничего, кроме того, что "здесь есть место, куда можно вставить слово".

Тем не менее, если когда-нибудь появится потребность в создании облака слов, то стоит приглядеться к осевым линиям и возможности при помощи них что-то выразить. Например, представим, что для каждого слова из некоторой совокупности популярных терминов науки о данных имеется два числа между 0 и 100: первое обозначает частотность его появлений в объявлениях о вакансиях, второе — частотность его появлений в резюме:

```
data = [{"big data",100,15}, {"Hadoop",95,25}, {"Python",75,50},
        {"R",50,40}, {"machine learning",80,20}, {"statistics",20,60},
        {"data science",60,70}, {"analytics",90,3},
        {"team player",85,85}, {"dynamic",2,90}, {"synergies",70,0},
        {"actionable insights",40,30}, {"think out of the box",45,10},
        {"self-starter",30,50}, {"customer focus",65,15},
        {"thought leadership",35,35}]
```

Метод визуализации на основе облака слов сводится к размещению на странице слов, набранных броским шрифтом (рис. 20.1).

<sup>1</sup> Реплика из комедии "Бесплодные усилия любви" (середина 1590-х гг.) Уильяма Шекспира (1564–1616) — знаменитого английского поэта и драматурга. — *Прим. пер.*



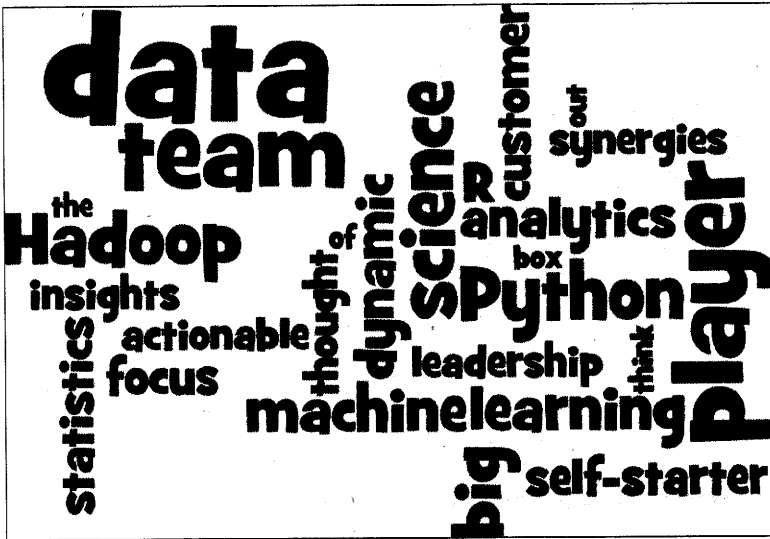


Рис. 20.1. Облако модных терминов

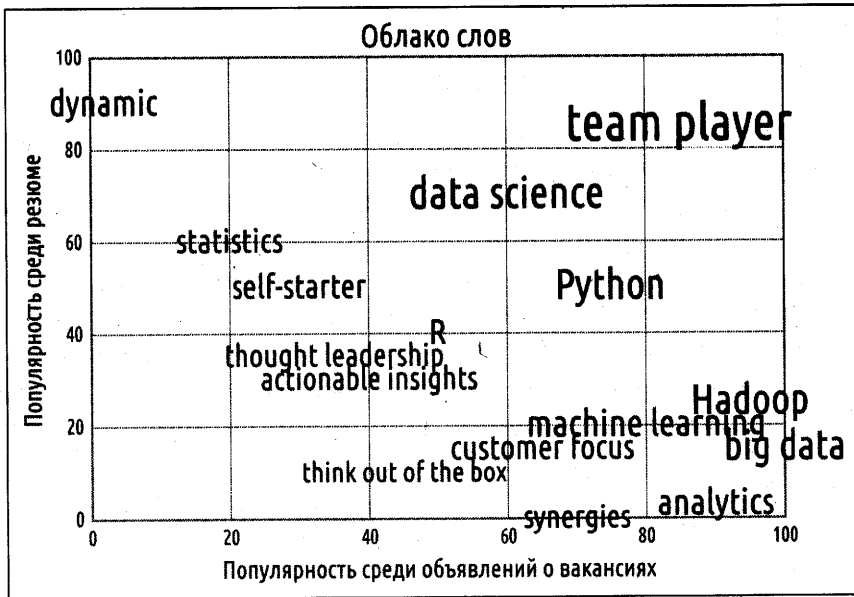


Рис. 20.2. Более осмысленное (и менее привлекательное) облако слов

Смотрится складно, но не говорит ни о чем. Более интересный подход может заключаться в том, чтобы разбросать слова так, что горизонтальное положение будет указывать на популярность среди объявлений, а вертикальное — на популярность среди резюме, что позволит получить визуализацию, которая будет нести в себе уже больше смысла (рис. 20.2):

```
# размер текста
def text_size(total):
```

```
"""равно 8, если total = 0, и 28, если total = 200"""
```

```
return 8 + total / 200 * 20
```

```
for word, job_popularity, resume_popularity in data:
    plt.text(job_popularity, resume_popularity, word,
             ha='center', va='center',
             size=text_size(job_popularity + resume_popularity))
plt.xlabel("Популярность среди объявлений о вакансиях")
plt.ylabel("Популярность среди резюме")
plt.axis([0, 100, 0, 100])
plt.xticks([])
plt.yticks([])
plt.show()
```

## N-граммные модели языка

Директор по поисковому маркетингу хочет создать тысячи веб-страниц, посвященных науке о данных, благодаря чему сайт DataSciencester станет ранжироваться выше среди результатов поиска терминов, связанных с наукой о данных. (Вы пытаетесь объяснить, что алгоритмы поисковых систем достаточно умны, и что это на самом деле не сработает, но он отказывается слушать.)

Естественно, сам писать тысячи веб-страниц он не будет, но и платить орде "контент-стратегов", чтобы те сделали это, тоже не хочет. Напротив, он интересуется у вас, можно ли создать эти веб-страницы программно. Для этого понадобится каким-то образом смоделировать язык.

Один из приемов заключается в том, чтобы вычислить статистическую модель языка, начав с корпуса документов. В нашем случае мы начнем с очерка Майка Лукидеса "Что такое наука о данных?"<sup>2</sup>

Как и в *главе 9*, воспользуемся библиотеками `requests` и `BeautifulSoup` для извлечения данных. Есть, правда, пара вопросов, на которые стоит обратить внимание.

Во-первых, апострофы в тексте на самом деле представлены символом Юникода `u"\u2019"`. Создадим вспомогательную функцию, которая будет менять их на нормальные апострофы:

```
def fix_unicode(text):
    return text.replace(u"\u2019", "'")
```

Во-вторых, получив текст веб-страницы, необходимо разбить его на последовательность слов и точек (благодаря чему можно распознать конец предложений). Это можно сделать при помощи метода `re.findall()`:

```
from bs4 import BeautifulSoup
import requests
```

<sup>2</sup> <https://www.oreilly.com/ideas/what-is-data-science>.

```

url = "http://radar.oreilly.com/2010/06/what-is-data-science.html"
html = requests.get(url).text
soup = BeautifulSoup(html, 'html5lib')

content = soup.find("div", "entry-content") # найти div с классом entry-content
regex = r"[\w']+|[\.]" # совпадает со словом или точкой

document = []

for paragraph in content("p"):
    words = re.findall(regex, fix_unicode(paragraph.text))
    document.extend(words)

```

Разумеется, можно (и, вероятно, следует) продолжить очистку данных. В документе еще осталось некоторое количество постороннего текста (первое слово "Section"), несколько разбросанных повсюду заголовков и списков и некорректное разбиение текста по точке в середине фразы ("Web 2.0"). Сказав это, все же продолжим работать с документом в том виде, в котором он нам достался.

Теперь, при наличии текста в виде последовательности слов, можно смоделировать язык следующим образом: имея некое исходное слово (скажем, "book"), обратимся ко всем словам, которые в исходных документах следуют за ним (здесь "isn't", "a", "shows", "demonstrates" и "teaches"). Произвольно выберем одно из них в качестве следующего слова и повторим итерационный процесс до тех пор, пока не получим точку, которая означает конец предложения. Такая процедура называется *биграммной языковой моделью*, поскольку она полностью обуславливается частотностями биграмм (пар слов) в исходных данных.

Как быть с исходным словом? Его можно выбрать произвольно из тех слов, которые *следуют* за точкой. Сначала предварительно рассчитаем возможные переходы от слова к слову. Напомним, что функция `zip` останавливается, как только любой из ее входящих списков заканчивается, благодаря чему выражение `zip(document, document[1:])` даст пары элементов, которые в документе `document` следуют строго один за другим:

```

bigrams = zip(document, document[1:]) # биграммы
transitions = defaultdict(list) # переходы
for prev, current in bigrams:
    transitions[prev].append(current)

```

Теперь все готово для генерирования предложений:

```

# сгенерировать при помощи биграмм
def generate_using_bigrams():
    current = "." # означает, что следующее слово начинает предложение
    result = []
    while True:
        next_word_candidates = transitions[current] # биграммы (current, _)
        current = random.choice(next_word_candidates) # выбрать произвольно
        result.append(current) # добавить к результатам
        if current == ".": return " ".join(result) # если ".", то закончить

```

Сгенерированные предложения выглядят полной тарабарщиной, но такой, что если разместить их на сайте, то они вполне могут сойти за глубокомысленные высказывания на тему науки о данных. Например:

If you may know which are you want to data sort the data feeds web friend someone on trending topics as the data in Hadoop is the data science requires a book demonstrates why visualizations are but we do massive correlations across many commercial disk drives in Python language and creates more tractable form making connections then use and uses it to solve a data.

Биграммная языковая модель.

Предложения можно сделать менее бредовыми, если обратиться к *триграммам*, т. е. тройкам последовательных слов. (В более общем плане можно обратиться к *n-граммам*, состоящим из *n* последовательных слов, однако трех для примера будет достаточно.) Теперь переходы будут зависеть от *двух* предыдущих слов:

```
trigrams = zip(document, document[1:], document[2:])
trigram_transitions = defaultdict(list) # триграммные переходы
starts = []
```

```
for prev, current, next in trigrams:
```

```
    if prev == ".": # если предыдущее "слово" - точка,
        starts.append(current) # то это слово исходное
```

```
    trigram_transitions[(prev, current)].append(next)
```

Следует учесть, что теперь требуется отслеживать исходные слова отдельно. Генерирование предложений выполняется практически в том же виде:

```
# сгенерировать при помощи триграмм
def generate_using_trigrams():
    current = random.choice(starts) # выбрать произвольное исходное слово
    prev = "." # и поставить перед ним '.'
    result = [current]
    while True:
        next_word_candidates = trigram_transitions[(prev, current)]
        next_word = random.choice(next_word_candidates)
        prev, current = current, next_word
        result.append(current)
    if current == ".":
        return " ".join(result)
```

В результате получим более осмысленные предложения, например:

In hindsight MapReduce seems like an epidemic and if so does that give us new insights into how economies work That's not a question we could even have asked a few years there has been instrumented.

Триграммная языковая модель.

Естественно, они выглядят более осмысленными, потому что на каждом шаге у процесса генерации остается меньше выбора слов, а для многих шагов — всего один. Следовательно, модель часто будет генерировать предложения (или по крайней мере длинные словосочетания), которые существуют в исходных данных дословно. Ситуацию можно исправить, если предоставить больше данных; помимо этого, модель будет работать лучше, если собрать  $n$ -граммы из нескольких очерков, посвященных науке о данных.

## Грамматики

Другой способ моделирования языка основан на *грамматиках*, т. е. правилах генерации допустимых предложений. Из начальной школы мы знаем о частях речи и о том, как они сочетаются. Например, если у вас был отвратительный учитель родного языка, то вы могли бы утверждать, что предложение обязательно должно состоять из *существительного*, за которым следует *глагол*. Тогда при наличии списка существительных и глаголов данное правило позволит генерировать предложения.

Определим грамматику чуть посложнее:

```
grammar = {
    "_S" : ["_NP _VP"],
    "_NP" : ["_N",
            "_A _NP _P _A _N"],
    "_VP" : ["_V",
            "_V _NP"],
    "_N" : ["data science", "Python", "regression"],
    "_A" : ["big", "linear", "logistic"],
    "_P" : ["about", "near"],
    "_V" : ["learns", "trains", "tests", "is"]
}
```

Условимся, что каждая запись в приведенной грамматике называется правилом подстановки, где имена, начинающиеся с символа подчеркивания, обозначают *нетерминальные лексемы* (или правила), которые нуждаются в дальнейшем расширении, а остальные имена — это *терминалы*, которые в дальнейшей обработке не нуждаются.

Так, например, "\_S" — это нетерминальная лексема, описывающая правило подстановки для "предложения", которое порождает "\_NP" ("группу существительного"), за которой следует "\_VP" ("группа глагола").

Нетерминальная лексема группы глагола "\_VP" может порождать нетерминальную лексему глагола "\_V" ("глагол") либо нетерминальную лексему глагола, за которой идет нетерминальная лексема группы существительного.

Обратите внимание, что нетерминальная лексема группы существительного "\_NP" содержит саму себя в правой части одного из своих правил подстановки. Грамматики могут быть рекурсивными, что позволяет даже таким грамматикам с конеч-

ным числом состояний, как эта, генерировать бесконечное число разных предложений.

Каким образом на основе такой грамматики генерируют предложения? Начнем со списка, содержащего лексему предложения "\_s", и затем будем неоднократно расширять каждую лексему, подставляя вместо нее правую часть одного из правил подстановки, выбираемого случайно. Обработка прекращается при наличии списка, состоящего исключительно из терминалов.

Например, одна из таких последовательностей может выглядеть следующим образом:

```
['_s']
['_NP', '_VP']
['_N', '_VP']
['Python', '_VP']
['Python', '_V', '_NP']
['Python', 'trains', '_NP']
['Python', 'trains', '_A', '_NP', '_P', '_A', '_N']
['Python', 'trains', 'logistic', '_NP', '_P', '_A', '_N']
['Python', 'trains', 'logistic', '_N', '_P', '_A', '_N']
['Python', 'trains', 'logistic', 'data science', '_P', '_A', '_N']
['Python', 'trains', 'logistic', 'data science', 'about', '_A', '_N']
['Python', 'trains', 'logistic', 'data science', 'about', 'logistic', '_N']
['Python', 'trains', 'logistic', 'data science', 'about', 'logistic', 'Python']
```

Как это реализовать? Для начала создадим простую вспомогательную функцию для идентификации терминалов:

```
# терминал?
def is_terminal(token):
    return token[0] != "_"
```

Далее нужно написать функцию для преобразования списка лексем в предложение. Будем искать первую нетерминальную лексему. Если такая отсутствует, то это значит, что имеем законченное предложение, и обработка прекращается.

Если же нетерминальная лексема имеется, то случайным образом выбираем для нее одно из правил подстановки. Если в результате подстановки получится терминал (т. е. слово), то просто подставим его вместо лексемы. В противном случае, имеем дело с последовательностью разделенных пробелами нетерминальных лексем, которые надо преобразовать (методом `split()`) в список и затем присоединить к существующим лексемам. В любом случае, повторим итерационный процесс с новым набором лексем.

Собрав все вместе, получим:

```
# расширить лексемы (применить правила подстановки) при заданной
# грамматике
def expand(grammar, tokens):
    for i, token in enumerate(tokens):
```

```

# пропустить терминалы
if is_terminal(token): continue

# если мы тут, значит, нашли нетерминальную лексему,
# произвольно выбрать для нее подстановку
replacement = random.choice(grammar[token])

if is_terminal(replacement):
    tokens[i] = replacement
else:
    tokens = tokens[:i] + replacement.split() + tokens[(i+1):]

# теперь вызвать функцию expand с новым списком лексем
return expand(grammar, tokens)

# если мы тут, значит, имеем одни терминалы; закончить обработку
return tokens

```

И теперь можно начинать генерировать предложения:

```

# сгенерировать предложение с заданной грамматикой
def generate_sentence(grammar):
    return expand(grammar, ["_S"])

```

Попробуйте изменить грамматику — добавить больше слов, больше правил, другие части речи — пока не будете готовы создать для вашей компании столько страниц текста, сколько нужно.

Граматики, как ни странно, дают более интересные результаты, если их использовать в противоположном направлении. К примеру, можно применить грамматику к заданному предложению с целью выполнить его синтаксический разбор, чтобы потом определить подлежащие и сказуемые и установить его смысл.

Уметь применять науку о данных для порождения текста — это, конечно, классная вещь; однако ее применение для того, чтобы *понять* смысл текста, сродни волшебству. (См. разд. "Для дальнейшего изучения" главы 16, где указаны библиотеки, которыми можно воспользоваться для этих целей.)

## Ремарка: метод сэмплирования по Гиббсу

Генерировать выборки из некоторых распределений достаточно просто. Следующие выражения дадут равномерно распределенные случайные величины:

```
random.random()
```

и нормально распределенные случайные величины:

```
inverse_normal_cdf(random.random())
```

Однако извлечение выборок из некоторых распределений вызывает затруднение. *Сэмплирование по Гиббсу* (Gibbs sampling) — это метод извлечения выборок из

многомерных распределений, когда известны лишь некоторые условные распределения.

Например, представим бросание двух игральных кубиков. Пусть  $x$  будет значением первого кубика, а  $y$  — сумма кубиков. Допустим, нужно сгенерировать большое число пар  $(x, y)$ . В этом случае легко извлечь выборки непосредственно:

```
# бросить кубик
def roll_a_die():
    return random.choice([1, 2, 3, 4, 5, 6])

# прямая выборка
def direct_sample():
    d1 = roll_a_die()
    d2 = roll_a_die()
    return d1, d1 + d2
```

Но, предположим, что известны только условные распределения. Распределение  $y$  в зависимости от  $x$  получить легко — если известно значение  $x$ , то  $y$  равномерно будет равно  $x + 1$ ,  $x + 2$ ,  $x + 3$ ,  $x + 4$ ,  $x + 5$  или  $x + 6$ :

```
# случайное y в зависимости от x
def random_y_given_x(x):
    """равномерное значение будет x + 1, x + 2, ..., x + 6"""
    return x + roll_a_die()
```

В обратную сторону уже сложнее. Например, если известно, что  $y = 2$ , тогда неизбежно  $x = 1$  (поскольку единственный случай, когда сумма граней двух кубиков равна 2, только если обе равны 1). Если известно, что  $y = 3$ , то  $x$  равномерно равно 1 или 2. Аналогичным образом, если  $y = 11$ , то  $x$  должно равняться 5 или 6:

```
# случайное x в зависимости от y
def random_x_given_y(y):
    if y <= 7:
        # если сумма <= 7, первый кубик равномерно будет равен
        # 1, 2, ..., (сумма - 1)
        return random.randrange(1, y)
    else:
        # если сумма > 7, первый кубик равномерно будет равен
        # (сумма - 6), (сумма - 5), ..., 6
        return random.randrange(y - 6, 7)
```

Метод сэмплирования по Гиббсу работает так: начинаем с любого (допустимого) значения  $x$  и  $y$  и затем многократно чередуем, подставляя вместо  $x$  случайное значение, выбранное в зависимости от  $y$ , и вместо  $y$  — случайное значение, выбранное в зависимости от  $x$ . После ряда итераций полученные значения  $x$  и  $y$  будут представлять собой выборку из безусловного совместного распределения:

```
# выборка по Гиббсу
def gibbs_sample(num_iters=100):
    x, y = 1, 2 # на самом деле не имеет значения, какие числа
```



```

for _ in range(num_iters):
    x = random_x_given_y(y) # случайное x в зависимости от y
    y = random_y_given_x(x) # случайное y в зависимости от x
return x, y

```

Можно проверить, что этот алгоритм даст результаты, аналогичные непосредственной выборке:

```

# сравнить распределения
def compare_distributions(num_samples=1000):
    counts = defaultdict(lambda: [0, 0])
    for _ in range(num_samples):
        counts[gibbs_sample()][0] += 1
        counts[direct_sample()][1] += 1
    return counts

```

Данный метод будет использован в следующем разделе.

## Тематическое моделирование

При создании рекомендательной системы "Аналитики, которых вы должны знать" в главе 1 отыскивались точные совпадения с темами, которыми пользователи интересуются.

Более продвинутый подход к пониманию сфер интересов пользователей заключается в попытке установить более общие *тематики*, которые лежат в основе интересующих пользователей тем. Для определения общих тематик в наборе документов обычно используется *метод латентного размещения Дирихле* (latent Dirichlet allocation). Давайте применим его к документам, состоящим из тем, интересующих каждого пользователя.

Данный метод имеет некоторую аналогию с наивным байесовским классификатором, разработанным в главе 13, в том, что в нем принимается вероятностная модель документов. Опуская более трудные математические подробности, отметим, что для целей настоящего изложения эта модель предполагает следующие аспекты.

- ◆ Имеется некоторое постоянное число  $K$  общих тематик.
- ◆ Имеется случайная величина, которая каждой тематике ставит в соответствие распределение вероятностей на множестве слов. Это распределение следует представить, как вероятность наблюдать слово  $w$  в заданной тематике  $k$ .
- ◆ Имеется еще одна случайная величина, которая каждому документу назначает распределение вероятностей на множестве тематик. Это распределение следует представить, как смесь тематик в документе  $d$ .
- ◆ Каждое слово в документе генерируется сначала путем случайного выбора тематики (из распределения тематик для данного документа) и затем случайного выбора слова (из распределения слов для данной тематики).

В частности, пусть имеется набор документов `documents`, где каждый документ — это список слов, а также соответствующий набор тематик документов `document_`

topics, из которого каждому слову в каждом документе назначается некая тематика (здесь число между 0 и  $K - 1$ ).

Благодаря этому пятое слово в четвертом документе равно:

```
documents[3][4]
```

а тематика, из которой это слово выбрано, равна:

```
document_topics[3][4]
```

Это весьма явным образом определяет распределение каждого документа на множестве тематик и неявным образом — распределение каждой тематики на множестве слов.

Можно оценить правдоподобие того, что тематика 1 служит причиной для появления определенного слова, сравнивая частотность, с которой тематика 1 производит это слово, с частотностью, с которой она производит *любое* слово. (При разработке спам-фильтра в *главе 13* аналогичным образом сравнивалась частотность появления каждого слова в спамном сообщении с суммарной частотностью слов в этом сообщении.)

Хотя эти тематики являются всего лишь числами, им можно дать описательные имена, обратившись к словам, которым присвоен самый тяжелый вес. Осталось лишь решить, как сгенерировать список тематик документов `document_topics`. Именно тут вступает в игру метод сэмплирования по Гиббсу.

Начнем с абсолютно случайного назначения тематик всем словам во всех документах. Пройдем по каждому документу слово за словом. По слову и документу сконструируем веса для каждой тематики, которые зависят от (текущего) распределения тематик в этом документе и (текущего) распределения слов для этой тематики. После этого применим эти веса для выборки новой тематики, соответствующей этому слову. Если повторить этот итерационный процесс многократно, то в итоге получим совместную выборку из распределения "тематика-слово" и распределения "документ-тематика".

Для начала потребуется функция, которая будет в случайном порядке выбирать индекс из произвольного множества весов:

```
# случайная выборка индекса из множества (в данном случае - весов)
def sample_from(weights):
    """возвращает i с вероятностью по формуле weights[i] / sum(weights)"""
    total = sum(weights)
    rnd = total * random() # равномерно между 0 и суммой
    for i, w in enumerate(weights):
        rnd -= w # вернуть наименьший i, такой что
        if rnd <= 0: return i # weights[0] + ... + weights[i] >= rnd
```

Например, если передать ей веса [1, 1, 3], то одну пятую случаев она будет возвращать нулевой индекс, одну пятую — 1 и три пятых — 2.

Документы представлены списками тем, которыми пользователи интересуются, в следующем виде:

```
documents = [
    ["Hadoop", "Big Data", "HBase", "Java", "Spark", "Storm", "Cassandra"],
    ["NoSQL", "MongoDB", "Cassandra", "HBase", "Postgres"],
    ["Python", "scikit-learn", "scipy", "numpy", "statsmodels", "pandas"],
    ["R", "Python", "statistics", "regression", "probability"],
    ["machine learning", "regression", "decision trees", "libsvm"],
    ["Python", "R", "Java", "C++", "Haskell", "programming languages"],
    ["statistics", "probability", "mathematics", "theory"],
    ["machine learning", "scikit-learn", "Mahout", "neural networks"],
    ["neural networks", "deep learning", "Big Data", "artificial intelligence"],
    ["Hadoop", "Java", "MapReduce", "Big Data"],
    ["statistics", "R", "statsmodels"],
    ["C++", "deep learning", "artificial intelligence", "probability"],
    ["pandas", "R", "Python"],
    ["databases", "HBase", "Postgres", "MySQL", "MongoDB"],
    ["libsvm", "regression", "support vector machines"]
]
```

Попробуем найти тематики при  $K = 4$ .

Для расчета весов выборки необходимо отслеживать несколько показателей. Сначала создадим для них структуры данных.

Сколько раз каждая тематика назначается каждому документу:

```
# список объектов Counter, один для каждого документа
document_topic_counts = [Counter() for _ in documents]
```

Сколько раз каждое слово назначается каждой тематике:

```
# список объектов Counter, один для каждой тематики
topic_word_counts = [Counter() for _ in range(K)]
```

Суммарное число слов, назначенное каждой тематике:

```
# список чисел, один для каждой тематики
topic_counts = [0 for _ in range(K)]
```

Суммарное число слов, содержащихся в каждом документе:

```
# список чисел, один для каждого документа
document_lengths = map(len, documents)
```

Число различных слов:

```
distinct_words = set(word for document in documents for word in document)
W = len(distinct_words)
```

И число документов:

```
D = len(documents)
```

Например, после того как они будут заполнены, можно найти, скажем, число слов в документе `documents[3]`, связанных с тематикой 1:

```
document_topic_counts[3][1]
```

Кроме этого, можно найти, сколько раз слово *nlp* ассоциируется с тематикой 2:

```
topic_word_counts[2]["nlp"]
```

Теперь все готово для определения функций условной вероятности. Как и в главе 13, каждая из них имеет сглаживающий член, который гарантирует, что каждая тематика будет иметь ненулевую вероятность быть выбранной в любом документе и что каждое слово будет иметь ненулевую вероятность быть выбранным в любой тематике:

```
# вероятность тематики в зависимости от документа
```

```
def p_topic_given_document(topic, d, alpha=0.1):
```

```
    """доля слов в документе _d_, которые
    назначаются тематике _topic_ (плюс некоторое сглаживание)"""
```

```
    return ((document_topic_counts[d][topic] + alpha) /
            (document_lengths[d] + K * alpha))
```

```
# вероятность слова в зависимости от тематики
```

```
def p_word_given_topic(word, topic, beta=0.1):
```

```
    """доля слов, назначаемых тематике _topic_, которые
    равны _word_ (плюс некоторое сглаживание)"""
```

```
    return ((topic_word_counts[topic][word] + beta) /
            (topic_counts[topic] + W * beta))
```

Эти функции используются для создания весов, необходимых для обновления тематик:

```
# вес тематики
```

```
def topic_weight(d, word, k):
```

```
    """при наличии документа и слова в этом документе,
    вернуть вес k-й темы"""
```

```
    return p_word_given_topic(word, k) * p_topic_given_document(k, d)
```

```
# выбрать новую тематику
```

```
def choose_new_topic(d, word):
```

```
    return sample_from([topic_weight(d, word, k)
                       for k in range(K)])
```

Имеются веские математические причины, почему функция `topic_weight` определена именно таким образом, однако их детализация увела бы далеко вглубь специализированной области. Надеюсь, по крайней мере, интуитивно понятно, что — при заданном слове и его документе — правдоподобие выбора любой тематики зависит одновременно от того, насколько правдоподобна эта тематика для документа и насколько правдоподобно это слово для этой тематики.

Вот и весь необходимый функциональный аппарат. Начнем с назначения каждого слова случайной тематике и заполнения показателей соответствующими значениями:

```

random.seed(0)
document_topics = [[random.randrange(K) for word in document]
                   for document in documents]

for d in range(D):
    for word, topic in zip(documents[d], document_topics[d]):
        document_topic_counts[d][topic] += 1
        topic_word_counts[topic][word] += 1
        topic_counts[topic] += 1

```

Задача состоит в том, чтобы получить совместные выборки из распределения "тематика-слово" и распределения "документ-тематика". Это достигается при помощи варианта метода сэмплирования по Гиббсу с использованием условных вероятностей, которые были определены ранее:

```

for iter in range(1000):
    for d in range(D):
        for i, (word, topic) in enumerate(zip(documents[d],
                                              document_topics[d])):
            # удалить это слово/тематику из показателей,
            # чтобы оно не влияло на веса
            document_topic_counts[d][topic] -= 1
            topic_word_counts[topic][word] -= 1
            topic_counts[topic] -= 1
            document_lengths[d] -= 1

            # выбрать новую тематику на основе весов
            new_topic = choose_new_topic(d, word)
            document_topics[d][i] = new_topic

            # а теперь снова увеличить показатели
            document_topic_counts[d][new_topic] += 1
            topic_word_counts[new_topic][word] += 1
            topic_counts[new_topic] += 1
            document_lengths[d] += 1

```

Что собой представляют тематикки? Это просто цифры 0, 1, 2 и 3. Если для них нужны имена, то это следует сделать самим. Обратимся к пяти наиболее тяжеловесным словам по каждой тематике (табл. 20.1):

```

for k, word_counts in enumerate(topic_word_counts):
    for word, count in word_counts.most_common():
        if count > 0: print(k, word, count)

```

На их основе тематикам можно присвоить имена:

```

topic_names = ["Big Data and programming languages",
               "Python and statistics",
               "databases",
               "machine learning"]

```

Таблица 20.1. Наиболее распространенные слова по тематике

Тема 0	Тема 1	Тема 2	Тема 3
Java	R	HBase	regression
Big Data	statistics	Postgres	libsvm
Hadoop	Python	MongoDB	scikit-learn
deep learning	probability	Cassandra	machine learning
artificial intelligence	pandas	NoSQL	neural networks

Затем можно проверить, каким образом модель присваивает тематики темам, которые интересуют каждого пользователя:

```
for document, topic_counts in zip(documents, document_topic_counts):
    print(document)
    for topic, count in topic_counts.most_common():
        if count > 0:
            print(topic_names[topic], count)
    print()
```

В результате получится:

```
['Hadoop', 'Big Data', 'HBase', 'Java', 'Spark', 'Storm', 'Cassandra']
Big Data and programming languages 4 databases 3
['NoSQL', 'MongoDB', 'Cassandra', 'HBase', 'Postgres']
databases 5
['Python', 'scikit-learn', 'scipy', 'numpy', 'statsmodels', 'pandas']
Python and statistics 5 machine learning 1
```

И так далее. Учитывая союзы "and", которые были нужны в некоторых названиях тематик, возможно, следует использовать другие тематики, хотя, весьма вероятно, просто не хватает данных, чтобы успешно их извлечь.

## Для дальнейшего изучения

- ◆ Естественно-языковой инструментарий (Natural Language Toolkit, NLTK; <http://www.nltk.org/>) — это популярная (и достаточно исчерпывающая) библиотека средств обработки естественного языка на Python. Библиотека располагает собственной книгой (<http://www.nltk.org/book/>), которая доступна для чтения в сети.
- ◆ gensim (<http://radimrehurek.com/gensim/>) — это питоновская библиотека для тематического моделирования, которая предлагает больше возможностей, чем показанная здесь модель, написанная с чистого листа.

# Анализ социальных сетей

Ваши связи со всем, что вас окружает, в буквальном смысле определяют вас.

Аарон О'Коннелл<sup>1</sup>

Многие увлекательные задачи, связанные с анализом данных, могут быть плодотворно осмыслены в терминах *сетей*, состоящих из *узлов* определенного типа и *дуг*, которые их соединяют.

Например, друзья в Facebook представляют собой узлы сети, в которой дуги — это дружеские отношения. Менее очевидным примером является непосредственно сама Всемирная паутина, где узлами являются веб-страницы, а дугами — гиперссылки с одной страницы на другую.

Дружба в Facebook взаимна: если я дружу с вами, то вы неизбежно дружите со мной. В этом случае говорят, что дуги *ненаправленные*. А вот гиперссылки — нет: мой веб-сайт ссылается на веб-сайт Белого дома [whitehouse.gov](http://whitehouse.gov), однако (по необъяснимым для меня причинам) [whitehouse.gov](http://whitehouse.gov) отказывается ссылаться на мой веб-сайт. Такие типы дуг называют *направленными*. Мы рассмотрим оба вида сетей.

## Центральность по посредничеству

В *главе 1* были вычислены ключевые звенья в сети DataSciencester путем подсчета числа друзей каждого пользователя. Теперь, имеющегося в распоряжении функционального аппарата достаточно, чтобы обратиться к другим подходам. Напомним, что сеть (рис. 21.1) составляют пользователи:

```
users = [  
    { "id": 0, "name": "Hero" },  
    { "id": 1, "name": "Dunn" },  
    { "id": 2, "name": "Sue" },  
    { "id": 3, "name": "Chi" },  
    { "id": 4, "name": "Thor" },  
    { "id": 5, "name": "Clive" },  
    { "id": 6, "name": "Hicks" },  
    { "id": 7, "name": "Devin" },  
    { "id": 8, "name": "Kate" },  
    { "id": 9, "name": "Klein" }  
]
```

<sup>1</sup> Аарон Д. О'Коннелл (1981) — американский ученый в области экспериментальной квантовой физики, создатель первой в мире квантовой машины. — *Прим. пер.*

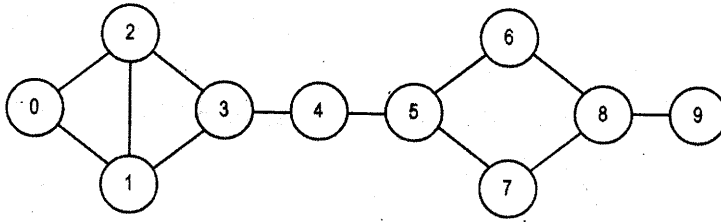


Рис. 21.1. Социальная сеть DataSciencester

и дружеские отношения:

```
friendships = [(0, 1), (0, 2), (1, 2), (1, 3), (2, 3), (3, 4),
              (4, 5), (5, 6), (5, 7), (6, 8), (7, 8), (8, 9)]
```

В питоновский словарь `dict` каждого пользователя также добавлен список друзей:

```
for user in users:
    user["friends"] = []

for i, j in friendships:
    # это работает, потому что users[i] - это пользователь, чей id равен i
    users[i]["friends"].append(users[j]) # добавить j как друга для i
    users[j]["friends"].append(users[i]) # добавить i как друга для j
```

Рассмотрение данной темы было отложено с чувством сомнения относительно понятия *центральности по степени узлов*, которое практически шло вразрез с интуитивным пониманием о том, кто является ключевым звеном сети.

Альтернативной метрикой является *центральность по посредничеству* (betweenness centrality), выявляющая людей, которые часто расположены на кратчайшем пути между парами других людей. В частности, центральность по посредничеству для узла  $i$  вычисляется путем сложения долей кратчайших путей между любыми другими парами вершин  $j$  и  $k$ , которые проходят через  $i$ .

Другими словами, чтобы выяснить центральность по посредничеству для Thor, нужно вычислить все кратчайшие пути между всеми парами людей за исключением Thor. И затем подсчитать количество кратчайших путей, которые проходят через пользователя Thor. Например, через Thor проходит единственный кратчайший путь между Chi (id 3) и Clive (id 5), и ни один из двух кратчайших путей между Hero (id 0) и Chi (id 3).

Итак, в качестве первого шага необходимо выяснить кратчайшие пути между всеми парами людей. Эту задачу эффективно решают несколько довольно сложных алгоритмов, однако (как и почти везде) воспользуемся алгоритмом, который менее эффективен, но легче для понимания.

Этот алгоритм (реализация поиска в ширину в графе) является одним из наиболее сложных в книге, так что остановимся на нем подробнее:

1. Нашей задачей является функция, которая получает пользователя `from_user` и находит все кратчайшие пути от него до остальных пользователей.



2. Представим путь в виде списка идентификаторов пользователей. Поскольку каждый путь начинается с `from_user`, его идентификатор включаться в список не будет. Вследствие этого длина списка, представляющего путь, будет равна длине самого пути.
3. Создадим словарь кратчайших путей до пункта назначения `shortest_paths_to`, в котором ключами будут идентификаторы пользователей, а значениями — списки путей, которые заканчиваются на пользователе с данным идентификатором ID. Если кратчайший путь всего один, то список будет содержать лишь его. Если же их несколько, то список будет содержать их все.
4. Кроме этого, создадим очередь `frontier`, содержащую пользователей, которых требуется исследовать в порядке их поступления. Они будут храниться в виде пар, состоящих из предыдущего и текущего пользователей (`prev_user, user`), благодаря которым мы узнаем, каким образом мы добрались до каждого из них. Очередь будет инициализирована всеми соседями пользователя `from_user`. (Несколько слов об очередях, которые упоминаются тут впервые. Это структуры данных, оптимизированные для операций добавления элементов в конец последовательности и удаления из ее начала. В языке Python они реализованы в виде двусвязной очереди `collections.deque`.)
5. Во время исследования графа, когда обнаруживаются новые соседи, кратчайшие пути до которых еще неизвестны, они будут добавляться в конец очереди с тем, чтобы исследовать их позже. При этом переменной `prev_user`, обозначающей предыдущего пользователя, будет присвоено значение текущего пользователя.
6. Когда из очереди изымается пользователь, который никогда раньше не встречался, это означает, что, безусловно, к нему найден один или несколько кратчайших путей — каждый из кратчайших путей к `prev_user` с добавлением одного лишнего шага.
7. Когда из очереди изымается пользователь, который уже встречался, то либо найден еще один кратчайший путь (в этом случае его следует добавить), либо найден более длинный путь (и в этом случае он не добавляется).
8. Когда в очереди больше нет пользователей, то весь граф (или, по крайней мере, та его часть, которая достижима из исходного пользователя) был исследован, и работа завершается.

Все это можно собрать в одну (большую) функцию:

```
from collections import deque
```

```
# кратчайшие пути от пользователя from_user
```

```
def shortest_paths_from(from_user):
```

```
    # словарь из "user_id" до ВСЕХ кратчайших путей к этому пользователю
```

```
    shortest_paths_to = { from_user["id"] : [[]] }
```

```
    # двусвязная очередь (предыдущий пользователь, следующий пользователь),
```

```
    # с которой нужно сверяться; инициализируется всеми парами,
```

```

# состоящими из исходного пользователя
# и его друзьями (from_user, friend_of_from_user)
frontier = deque((from_user, friend)
                 for friend in from_user["friends"])

# продолжать, пока очередь не пуста
while frontier:
    prev_user, user = frontier.popleft() # удалить пользователя,
    user_id = user["id"] # который стоит в очереди первым

    # в силу особенности добавления элементов к очереди с неизбежностью
    # некоторые из кратчайших путей к prev_user уже известны
    paths_to_prev_user = shortest_paths_to[prev_user["id"]]
    new_paths_to_user = [path + [user_id] for path in paths_to_prev_user]

    # возможно, кратчайший путь уже известен
    old_paths_to_user = shortest_paths_to.get(user_id, [])

    # каков кратчайший путь до этого места, которое уже встречалось ранее?
    if old_paths_to_user:
        min_path_length = len(old_paths_to_user[0])
    else:
        min_path_length = float('inf')

    # хранить только пути, которые не слишком длинные
    # и действительно новые
    new_paths_to_user = [path
                        for path in new_paths_to_user
                        if len(path) <= min_path_length
                        and path not in old_paths_to_user]

    shortest_paths_to[user_id] = old_paths_to_user + new_paths_to_user

    # добавить к очереди frontier не встречавшихся ранее соседей
    frontier.extend((user, friend)
                   for friend in user["friends"]
                   if friend["id"] not in shortest_paths_to)

return shortest_paths_to

```

Теперь в каждом узле сети можно разместить словари dict с кратчайшими путями:

```

# свойство shortest_paths содержит кратчайшие пути для пользователя user
for user in users:
    user["shortest_paths"] = shortest_paths_from(user)

```

Наконец все готово для вычисления центральности по посредничеству. Поскольку для каждой пары узлов  $i$  и  $j$  известны  $n$  кратчайших путей из  $i$  в  $j$ , то для каждого из этих путей добавим  $1/n$  к центральности каждого узла на этом пути:

```
# свойство betweenness centrality содержит значение
# центральности по посредничеству для пользователя user
for user in users:
    user["betweenness centrality"] = 0.0

for source in users:
    source_id = source["id"]
    for target_id, paths in source["shortest_paths"].iteritems():
        if source_id < target_id: # чтобы избежать дублирования
            num_paths = len(paths) # сколько кратчайших путей?
            contrib = 1 / num_paths # вклад в центральность
            for path in paths:
                for id in path:
                    if id not in [source_id, target_id]:
                        users[id]["betweenness centrality"] += contrib
```

Как показано на рис. 21.2, пользователи 0 и 9 имеют нулевую центральность (т. к. ни один не находится на кратчайшем пути между другими пользователями), тогда как пользователи 3, 4 и 5 имеют высокие центральности (т. к. все трое находятся на нескольких кратчайших путях).

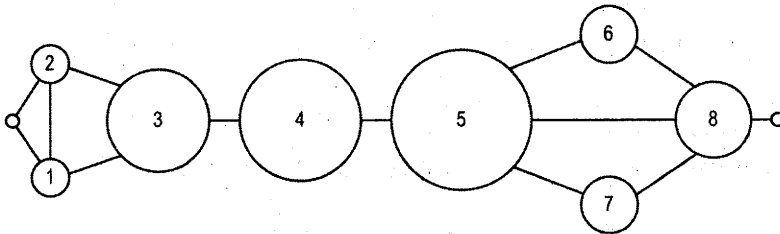


Рис. 21.2. Сеть после измерения метрикой центральности по посредничеству



В целом показатели центральности сами по себе не имеют большого значения. Практическую значимость имеет то, как показатели для каждого узла соотносятся с показателями для других узлов.

Еще одна мера, к которой можно обратиться, называется *центральностью по близости* (closeness centrality). Сначала для каждого пользователя вычисляется его *удаленность*, т. е. сумма длин его кратчайших путей к любому другому пользователю. Так как кратчайшие пути между каждой парой узлов уже вычислены, то сложить их длины не представляет труда. (Если кратчайших путей несколько, то все они будут иметь одинаковые длины, поэтому можно просто взять первый из них.)

```
# удаленность
def farness(user):
    """сумма длин кратчайших путей к любому другому пользователю"""
    return sum(len(paths[0])
               for paths in user["shortest_paths"].values())
```

И теперь легко вычисляется центральность по близости (рис. 21.3):

```
# свойство closeness_centrality содержит значение
# центральности по близости для пользователя user
for user in users:
    user["closeness_centrality"] = 1 / farness(user)
```

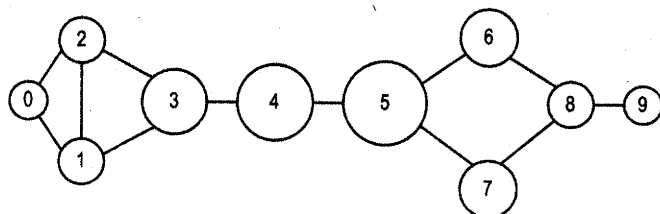


Рис. 21.3. Сеть после измерения метрикой центральности по близости

Здесь гораздо меньше изменений — ближайшие к центру узлы по-прежнему лежат довольно далеко от узлов на периферии.

Как можно было убедиться, вычисление кратчайших путей создает некоторые проблемы. По этой причине метрики центральности по посредничеству и центральности по близости используются в больших сетях не так часто. Чаще используют менее интуитивно понятную (но, как правило, простую для вычисления) метрику *центральности собственного вектора*.

## Центральность собственного вектора

Чтобы начать обсуждать *центральность собственного вектора* (eigenvector centrality), сначала следует поговорить о самих собственных векторах, а чтобы говорить о них, следует остановиться на умножении матриц.

### Умножение матриц

Если  $\mathbf{A}$  — это матрица размера  $n_1 \times k_1$ , а  $\mathbf{B}$  — матрица размера  $n_2 \times k_2$  и, если  $k_1 = n_2$ , тогда их произведение  $\mathbf{AB}$  есть матрица размера  $n_1 \times k_2$ , чья  $(i, j)$ -я запись имеет вид:

$$a_{i1}b_{1j} + a_{i2}b_{2j} + \dots + a_{ik}b_{kj},$$

что в точности соответствует скалярному произведению dot  $i$ -й строки матрицы  $\mathbf{A}$  (представленной как вектор) на  $j$ -й столбец матрицы  $\mathbf{B}$  (тоже представленной как вектор):

```
# запись в умноженной матрице
def matrix_product_entry(A, B, i, j):
    return dot(get_row(A, i), get_column(B, j))
```

После чего получим:

```
# умножение матриц
def matrix_multiply(A, B):
    n1, k1 = shape(A)
    n2, k2 = shape(B)
    if k1 != n2:
        raise ArithmeticError("несовместимые формы матриц!")

    return make_matrix(n1, k2, partial(matrix_product_entry, A, B))
```

Следует обратить внимание, что если **A** — это матрица размера  $n \times k$ , а **B** — матрица размера  $k \times 1$ , тогда **AB** — это матрица размером  $n \times 1$ . Если рассматривать вектор как матрицу из одного столбца, то можно представить **A** как функцию, которая отображает  $k$ -мерные векторы в  $n$ -мерные векторы, где функция — это именно умножение матриц.

Ранее векторы представлялись списками, что не совсем то же самое:

```
v = [1, 2, 3]          # v в виде списка
v_as_matrix = [[1],   # v в виде матрицы
               [2],
               [3]]
```

Поэтому понадобится несколько вспомогательных функций для прямого и обратного преобразования между двумя представлениями:

```
# преобразовать вектор в матрицу
def vector_as_matrix(v):
    """возвращает n x 1 матричное представление
    для вектора v (представленного списком)"""
    return [[v_i] for v_i in v]

# преобразовать матрицу в вектор
def vector_from_matrix(v_as_matrix):
    """возвращает векторное представление в виде списка значений
    для n x 1 матрицы"""
    return [row[0] for row in v_as_matrix]
```

Теперь при помощи функции умножения матриц `matrix_multiply` можно определить матричный оператор `matrix_operate`:

```
# матричный оператор (для преобразований и умножения)
def matrix_operate(A, v):
    v_as_matrix = vector_as_matrix(v)
    product = matrix_multiply(A, v_as_matrix)
    return vector_from_matrix(product)
```

Если **A** — это *квадратная* матрица, то данная операция отображает  $n$ -мерные векторы в другие  $n$ -мерные векторы. Может оказаться, что для некоторой матрицы **A** и вектора **v**, когда **A** действует на **v**, получается скалярное кратное для **v**. Иными словами, в результате получен вектор, который указывает в том же самом направ-

лении, что и  $v$ . Когда это происходит (и когда вдобавок  $v$  — это не вектор нулей), то  $v$  называется *собственным вектором*  $A$ , а множитель — *собственным числом*.

Один из возможных способов найти собственный вектор для матрицы  $A$  состоит в выборе начального вектора  $v$ , применении матричного оператора `matrix_operate`, нормировании результата и повторении до тех пор, пока итерационный процесс не достигнет схождения:

```
# найти собственный вектор
def find_eigenvector(A, tolerance=0.00001):
    guess = [random.random() for __ in A] # начальное приближение

    while True:
        result = matrix_operate(A, guess)
        length = magnitude(result)
        next_guess = scalar_multiply(1/length, result) # нормировать вектор

        # расстояние между текущим и новым приближениями
        # меньше критерия завершения итераций?
        if distance(guess, next_guess) < tolerance:
            return next_guess, length # собственный вектор, собственное число

    guess = next_guess # следующее приближение
```

По структуре возвращаемая переменная `guess` — это вектор, такой, что если применить к нему оператор `matrix_operate` и нормировать (т. е. чтобы он имел длину равную 1), то будет получен сам вектор (точнее вектор, очень близкий самому себе). И, следовательно, он является собственным вектором.

Не все матрицы вещественных чисел имеют собственные векторы и собственные числа. Например, матрица поворота:

```
rotate = [[ 0, 1],
          [-1, 0]]
```

поворачивает векторы на 90 градусов по часовой стрелке, и значит, вектор нулей будет единственным вектором, который она отображает в собственное скалярное кратное. Если попытаться вызвать `find_eigenvector(rotate)`, то функция войдет в бесконечный цикл. Более того, иногда могут заикливиться даже те матрицы, у которых собственные векторы существуют. Рассмотрим матрицу отражения:

```
flip = [[0, 1],
        [1, 0]]
```

Данная матрица отображает любой вектор  $[x, y]$  в  $[y, x]$ . Следовательно, например,  $[1, 1]$  — это собственный вектор с собственным числом, равным 1. Однако если начать со случайного вектора, у которого неравные координаты, то функция `find_eigenvector` будет бесконечно обмениваться координатами. (В библиотеках, написанных не с чистого листа, таких как NumPy, используются другие методы, которые в этом случае будут работать.) Тем не менее, тогда, когда `find_eigenvector`

возвращает результат, этот результат действительно является собственным вектором.

## Центральность

Как все это поможет разобраться в социальной сети DataSciencester?

Для начала необходимо представить связи в сети в виде матрицы смежности `adjacency_matrix`, чей  $(i, j)$ -й элемент равен либо 1 (если пользователи  $i$  и  $j$  — друзья), либо 0 (если нет):

```
# запись в матрице смежности
def entry_fn(i, j):
    return 1 if (i, j) in friendships or (j, i) in friendships else 0

n = len(users)
adjacency_matrix = make_matrix(n, n, entry_fn) # матрица смежности
```

Тогда центральность собственного вектора для каждого пользователя — это запись, соответствующая этому пользователю в собственном векторе, возвращаемом функцией `find_eigenvector` (рис. 21.4):

```
# получение центральностей собственного вектора
# на основе матрицы смежности
eigenvector_centralities, _ = find_eigenvector(adjacency_matrix)
```



По техническим причинам, которые выходят за рамки этой книги, договоримся, что любая ненулевая матрица смежности обязательно имеет собственный вектор всех неотрицательных значений. И к счастью, функция `find_eigenvector` находит его для матрицы смежности `adjacency_matrix`.

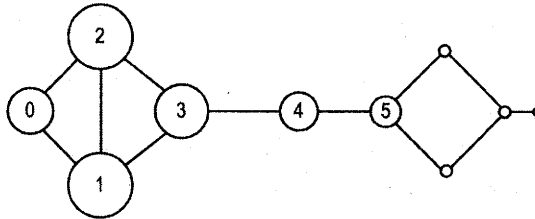


Рис. 21.4. Сеть после измерения метрикой центральности собственного вектора

Пользователями с высокой центральностью собственного вектора должны быть те, у кого много связей с людьми, которые сами имеют высокую центральность.

Здесь пользователи 1 и 2 находятся ближе всего к центру, т. к. у них обоих имеются три связи с людьми, которые сами находятся близко к центру. При удалении от них центральности людей неуклонно снижаются.

В такой небольшой сети, как эта, центральность собственного вектора ведет себя несколько хаотично. Если попробовать добавлять или удалять связи, то обнаружится, что небольшие изменения в сети могут кардинально изменить показатели центральности. Для более крупных сетей это нехарактерно.

До сих пор не было дано объяснений, почему собственный вектор приводит к разумному пониманию центральности. Собственный вектор означает, что если вычислить:

```
matrix_operate(adjacency_matrix, eigenvector_centralities)
```

то результатом будет скалярное кратное для центральностей собственного вектора `eigenvector_centralities`.

Если обратиться к тому, как работает умножение матриц, то оператор `matrix_operate` возвращает вектор, чей  $i$ -й элемент равен:

```
dot(get_row(adjacency_matrix, i), eigenvector_centralities)
```

который является в точности суммой центральностей собственных векторов пользователей, связанных с пользователем  $i$ .

Другими словами, центральности собственного вектора являются числами, одно на каждого пользователя, такими, что значение по каждому пользователю есть постоянный множитель суммы значений его соседей. В данном случае центральность означает наличие связи с людьми, которые сами центральны. Чем больше центральность, с которой вы непосредственно связаны, тем выше ваша собственная центральность. Такое определение, разумеется, имеет циклический характер — собственные векторы предоставляют способ вырваться из этой цикличности.

Еще один способ понять их — разобраться в том, что делает функция нахождения собственного вектора `find_eigenvector`. Ее работа начинается с назначения каждому узлу случайной центральности, после чего следующие два действия повторяются до тех пор, пока итерационный процесс не сойдется:

1. Дать каждому узлу значение новой центральности, которое равно сумме (старых) значений центральности своих соседей.
2. Нормировать вектор центральностей (привести к длине, равной 1).

Хотя лежащая в основе математика может сначала показаться несколько непонятной, сами расчеты довольно просты (в отличие, скажем, от центральности по посредничеству) и довольно легко выполняются даже на очень больших графах.

## Направленные графы и рейтинг PageRank

DataSciencester не получает достаточной поддержки, поэтому директор по доходам планирует перейти от модели дружеских отношений к модели *лидерства мнения*. Оказывается, никого особо не интересует, какие аналитики данных *дружат* друг с другом, однако агенты по найму научно-технического персонала очень интересуются тем, какие аналитики пользуются в своей среде заслуженным уважением.

В новой модели отслеживаются оценки в формате "источник — цель" (`source`, `target`), где представлены не взаимные связи, а связи, в которых пользователь `source` подтверждает, что пользователь `target` является потрясающим специалистом в области науки о данных (рис. 21.5). Эту асимметрию необходимо учесть:



```

# оценки
endorsements = [(0, 1), (1, 0), (0, 2), (2, 0), (1, 2),
                 (2, 1), (1, 3), (2, 3), (3, 4), (5, 4),
                 (5, 6), (7, 5), (6, 8), (8, 7), (8, 9)]

# свойство endorses содержит список лиц, которых пользователь user оценил;
# свойство endorsed_by содержит список лиц, которые оценили
# пользователя user
for user in users:
    user["endorses"] = [] # этот список отслеживает исходящие оценки
    user["endorsed_by"] = [] # этот список отслеживает поступающие оценки

for source_id, target_id in endorsements:
    users[source_id]["endorses"].append(users[target_id])
    users[target_id]["endorsed_by"].append(users[source_id])

```

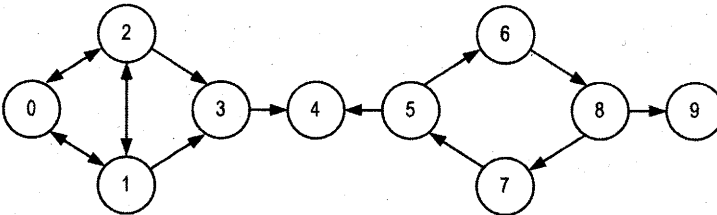


Рис. 21.5. Отношения авторитетности в сети DataSciencecenter

Теперь легко можно найти наиболее авторитетных `most_endorsed` аналитиков данных и продать эту информацию агентам по найму:

```

endorsements_by_id = [(user["id"], len(user["endorsed_by"]))
                      for user in users]

```

```

sorted(endorsements_by_id,
       key=lambda (user_id, num_endorsements): num_endorsements,
       reverse=True)

```

Тем не менее, "количество оценок" легко можно подделать. Для этого надо всего лишь завести фальшивые аккаунты, чтобы те поддержали вас, или договориться с друзьями поддержать друг друга (что пользователи 0, 1 и 2, кажется, сделали).

Усовершенствованная метрика могла бы учитывать *авторитетность* оценивающего. Оценки, полученные от людей, которые сами имеют много оценок, должны каким-то образом значить больше, чем оценки от людей лишь с несколькими оценками. В этом суть алгоритма PageRank, используемого в Google для ссылочного ранжирования веб-сайтов: в основе лежит "важность" сайтов, которые на них ссылаются, а у тех — "важность" других сайтов, которые ссылаются на них, и т. д.

(Если это несколько напоминает идею о центральности собственного вектора, то так и должно быть.)

Упрощенная версия алгоритма выглядит так:

1. В сети имеется суммарный рейтинг PageRank, равный 1.0 (или 100%).
2. Первоначально рейтинг распределен между узлами равномерно.
3. На каждом шаге значительная часть рейтинга каждого узла равномерно распределяется среди исходящих связей.
4. На каждом шаге оставшаяся часть рейтинга каждого узла равномерно распределяется среди всех узлов.

```
def page_rank(users, damping = 0.85, num_iters = 100):
```

```
    # первоначально распределить PageRank равномерно
    num_users = len(users)
    pr = { user["id"] : 1 / num_users for user in users }

    # это малая доля индекса PageRank,
    # которую каждый узел получает на каждой итерации
    base_pr = (1 - damping) / num_users

    for __ in range(num_iters):
        next_pr = { user["id"] : base_pr for user in users }
        for user in users:
            # распределить PageRank среди исходящих связей
            links_pr = pr[user["id"]] * damping
            for endorsee in user["endorses"]:
                next_pr[endorsee["id"]] += links_pr / len(user["endorses"])
        pr = next_pr

    return pr
```

Алгоритм PageRank (рис. 21.6) отождествляет пользователя id 4 (Thor) как аналитика данных с наивысшим рейтингом.

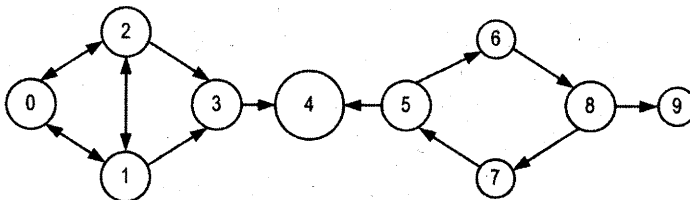


Рис. 21.6. Сеть после измерения метрикой PageRank

Несмотря на то, что у него меньше оценок (2), чем у пользователей 0, 1 и 2, его оценки несут в себе рейтинг от других оценок. Вдобавок, оба его сторонника поддержали только его, вследствие чего ему не нужно делиться их рейтингом с кем бы то ни было еще.

## Для дальнейшего изучения

- ◆ Помимо рассмотренных (надо признать, что рассмотренные метрические показатели центральности пользуются наибольшей популярностью) существует ряд других представлений о центральности (<https://en.wikipedia.org/wiki/Centrality>).
- ◆ NetworkX (<http://networkx.github.io/>) — это питоновская библиотека для выполнения сетевого анализа, которая располагает функциями для вычисления центральностей и визуализации графов.
- ◆ Gephi — это инструмент для визуализации сетей из разряда "от любви до ненависти" на основе графического интерфейса пользователя.

# Рекомендательные системы

О, природа, природа, и зачем же ты так нечестно поступаешь, когда посылаешь людей в этот мир с ложными рекомендациями!

Генри Филдинг<sup>1</sup>

Еще одна распространенная задача, связанная с аналитической обработкой данных, заключается в генерировании разного рода *рекомендаций*. Компания Netflix рекомендует фильмы, которые вы, возможно, захотите посмотреть, Amazon — товары, которые вы можете захотеть приобрести, Twitter — пользователей, за которыми вы могли бы последовать. В этой главе мы обратимся к нескольким способам использования данных для предоставления рекомендаций.

В частности, обратимся к набору данных `users_interests` о сферах интересов пользователей, который использовался ранее:

```
# темы, интересующие пользователей
users_interests = [
    ["Hadoop", "Big Data", "HBase", "Java", "Spark", "Storm", "Cassandra"],
    ["NoSQL", "MongoDB", "Cassandra", "HBase", "Postgres"],
    ["Python", "scikit-learn", "scipy", "numpy", "statsmodels", "pandas"],
    ["R", "Python", "statistics", "regression", "probability"],
    ["machine learning", "regression", "decision trees", "libsvm"],
    ["Python", "R", "Java", "C++", "Haskell", "programming languages"],
    ["statistics", "probability", "mathematics", "theory"],
    ["machine learning", "scikit-learn", "Mahout", "neural networks"],
    ["neural networks", "deep learning", "Big Data", "artificial intelligence"],
    ["Hadoop", "Java", "MapReduce", "Big Data"],
    ["statistics", "R", "statsmodels"],
    ["C++", "deep learning", "artificial intelligence", "probability"],
    ["pandas", "R", "Python"],
    ["databases", "HBase", "Postgres", "MySQL", "MongoDB"],
    ["libsvm", "regression", "support vector machines"]
]
```

При этом подумаем над вопросом генерирования для пользователя предложений относительно новых тем на основе тем, интересующих его в данное время.

<sup>1</sup> Генри Филдинг (1707–1754) — знаменитый английский писатель и драматург XVIII в., известен своим житейским юмором и сатирическим мастерством. Один из основоположников реалистического романа (см. [https://ru.wikipedia.org/wiki/Филдинг,\\_Генри](https://ru.wikipedia.org/wiki/Филдинг,_Генри)). — *Прим. пер.*

## Неавтоматическое кураторство

До появления Интернета за советом по поводу книг обычно ходили в библиотеку. Там библиотекарь был готов предложить книги, которые имеют отношение к сфере интересов читателя, либо аналогичные книги по желанию.

Учитывая ограниченное число пользователей DataSciencester и интересующих тем, нетрудно провести полдня, вручную раздавая рекомендации каждому пользователю по поводу новых тем. Но этот метод не очень хорошо поддается масштабированию, и ограничен личными познаниями и воображением, (Здесь нет ни капли намека на ограниченность чьих-то личных познаний и воображения.) Посмотрим, что можно сделать при помощи *данных*.

## Рекомендация популярных тем

Простейший подход заключается в рекомендации на основе популярности тем:

```
# популярные темы, интересующие пользователей
```

```
popular_interests = Counter(interest
                             for user_interests in users_interests
                             for interest in user_interests).most_common()
```

которые выглядят так:

```
[('Python', 4),
 ('R', 4),
 ('Java', 3),
 ('regression', 3),
 ('statistics', 3),
 ('probability', 3),
 # ...
]
```

После их вычисления может оказаться, что предложенные пользователю наиболее популярные темы его уже не интересуют:

```
# наиболее популярные новые темы
```

```
def most_popular_new_interests(user_interests, max_results=5):
    suggestions = [(interest, frequency)
                   for interest, frequency in popular_interests
                   if interest not in user_interests]
    return suggestions[:max_results]
```

Так, для пользователя 1, чья сфера интересов лежит в:

```
["NoSQL", "MongoDB", "Cassandra", "HBase", "Postgres"]
```

рекомендация пяти новых тем будет выглядеть следующим образом:

```
most_popular_new_interests(users_interests[1], 5)
```

```
# [('Python',4), ('R',4), ('Java',3), ('regression',3), ('statistics',3)]
```

Для пользователя 3, который уже интересуется многими из ниже следующих тем, будет получено:

```
[('Java', 3),
 ('HBase', 3),
 ('Big Data', 3),
 ('neural networks', 2),
 ('Hadoop', 2)]
```

Естественно, маркетинговый ход в стиле "многие интересуются языком Python, так, может, и вам тоже стоит" не самый убедительный. Если на сайт зашел новичок, о котором ничего неизвестно, то, возможно, это было бы лучшее, что можно сделать. Посмотрим, каким образом можно усовершенствовать рекомендации каждому пользователю, если основываться на интересующих его темах.

## Коллаборативная фильтрация на основе пользователя

Один из способов учесть сферу интересов пользователя заключается в том, чтобы отыскать пользователей, которые на него как-то *похожи*, и затем предложить темы, в которых эти пользователи заинтересованы.

Для этого потребуется метод измерения степени сходства двух пользователей. Здесь будет использована метрика, называемая *косинусным коэффициентом подобия* (cosine similarity). Пусть даны два вектора,  $v$  и  $w$ . Тогда косинусный коэффициент определяется следующим образом:

```
# косинусный коэффициент подобия
def cosine_similarity(v, w):
    return dot(v, w) / math.sqrt(dot(v, v) * dot(w, w))
```

Он измеряет "угол" между векторами  $v$  и  $w$ . Если векторы  $v$  и  $w$  — равнонаправлены, то числитель и знаменатель равны, и их косинусный коэффициент равен 1. Если они противоположно направленные, то их косинусный коэффициент равен  $-1$ . И если вектор  $v$  нулевой, а вектор  $w$  — нет (и наоборот), то результат  $\text{dot}(v, w)$  равен 0, и соответственно их косинусный коэффициент равен 0.

Применим его к векторам, состоящим из 0 и 1, где каждый вектор  $v$  обозначает темы, интересующие одного пользователя, т. е. его сферу интересов. Элемент  $v[i]$  равен 1, если пользователь указал  $i$ -ю тему, и 0 — в противном случае. Соответственно, под "сходством пользователей" будут пониматься "пользователи, чьи векторы интересующих их тем наиболее близко указывают в том же направлении". У людей с одинаковыми интересами коэффициент подобия будет равен 1, при отсутствии одинаковых интересов коэффициент подобия будет равен 0. В противном случае коэффициент подобия попадет в интервал между этими числами, где число, близкое к 1, будет указывать на большое "сходство", а число, близкое к 0, — на небольшое.

Неплохо начать с создания списка известных тем и (неявного) присваивания им индексов. Это можно сделать, воспользовавшись генератором последовательности

в виде множества, которое позволяет получить перечень тем без повторов, затем преобразовав полученный результат в список и упорядочив его. Первая тема в полученном списке будет с индексом 0 и т. д.:

```
# список интересных тем без повторов
unique_interests = sorted(list({ interest
                                for user_interests in users_interests
                                for interest in user_interests })))
```

В итоге получится список, который начинается со следующих записей:

```
['Big Data',
 'C++',
 'Cassandra', 'HBase',
 'Hadoop', 'Haskell',
 # ...
]
```

Далее для каждого пользователя следует создать вектор интересующих его тем из 0 и 1. Для этого надо лишь просмотреть список уникальных тем `unique_interests`, назначая 1, если пользователь заинтересован в теме, и 0 — если нет:

```
# создать вектор интересующих пользователя тем
def make_user_interest_vector(user_interests):
    """при заданном списке интересующих пользователя тем создать вектор,
    чей i-й элемент равен 1, если unique_interests[i] есть в списке,
    и 0 в противном случае"""
    return [1 if interest in user_interests else 0
            for interest in unique_interests]
```

Затем можно создать матрицу сфер интересов пользователей, отобразив эту функцию на каждый элемент списка списков интересующих тем (т. е. применив ее к каждому элементу списка при помощи функции `map`):

```
# матрица сфер интересов в формате (пользователь, интересующие темы),
# где список тем для каждого пользователя преобразован в 0 и 1
user_interest_matrix = map(make_user_interest_vector, users_interests)
```

Теперь элемент матрицы `user_interest_matrix[i][j]` равен 1, если пользователь *i* указал заинтересованность в теме *j*, и 0 — в противном случае.

Поскольку набор данных небольшой, вычислить попарные коэффициенты подобия между всеми пользователями не составляет труда:

```
# матрица сходств между пользователями
user_similarities = [[cosine_similarity(interest_vector_i, interest_vector_j)
                      for interest_vector_j in user_interest_matrix]
                     for interest_vector_i in user_interest_matrix]
```

Теперь элемент матрицы `user_similarities[i][j]` будет показывать коэффициент подобия между пользователями *i* и *j*.

Например, элемент матрицы `user_similarities[0][9]` равен 0.57, поскольку оба пользователя разделяют заинтересованность в Hadoop, Java и Big Data. С другой сторо-

ны, элемент матрицы `user_similarities[0][8]` равен всего 0.19, поскольку пользователи 0 и 8 имеют всего одну общую заинтересованность, а именно в Big Data.

В данном случае `user_similarities[i]` — это вектор коэффициентов подобия пользователя  $i$  с остальными пользователями. Этим вектором можно воспользоваться для написания функции, которая для заданного пользователя находит наиболее похожих пользователей. При этом, разумеется, ни сам пользователь, ни пользователи с нулевым коэффициентом не рассматриваются, а результаты сортируются по убыванию от наиболее до наименее похожих:

```
# пользователи наиболее похожие на пользователя user_id
def most_similar_users_to(user_id):
    pairs = [(other_user_id, similarity)           # найти других
             for other_user_id, similarity in     # пользователей с
               enumerate(user_similarities[user_id]) # ненулевым коэфф.
             if user_id != other_user_id and similarity > 0] # подобия

    return sorted(pairs,                          # отсортировать их
                  key=lambda (_, similarity): similarity, # сперва наиболее
                  reverse=True)                   # похожие
```

К примеру, если вызовем `most_similar_users_to(0)`, то получим:

```
[(9, 0.5669467095138409),
 (1, 0.3380617018914066),
 (8, 0.1889822365046136),
 (13, 0.1690308509457033),
 (5, 0.1543033499620919)]
```

Каким образом следует использовать этот результат для рекомендации пользователю новых тем? По каждой теме можно попросту просуммировать коэффициенты подобия других пользователей, которые интересуются данной темой:

```
# рекомендации на основе пользователя
def user_based_suggestions(user_id, include_current_interests=False):
    # суммировать все коэффициенты подобия
    suggestions = defaultdict(float) # рекомендации
    for other_user_id, similarity in most_similar_users_to(user_id):
        for interest in users_interests[other_user_id]:
            suggestions[interest] += similarity

    # преобразовать их в отсортированный список
    suggestions = sorted(suggestions.items(),
                        key=lambda (_, weight): weight,
                        reverse=True)

    # и (может быть) исключить уже имеющиеся темы
    if include_current_interests:
        return suggestions
```



```
else:
    return [(suggestion, weight)
            for suggestion, weight in suggestions
            if suggestion not in users_interests[user_id]]
```

Если вызвать `user_based_suggestions(0)`, то несколько первых рекомендуемых тем из всего списка будут следующими:

```
[('MapReduce', 0.5669467095138409),
 ('MongoDB', 0.50709255283711),
 ('Postgres', 0.50709255283711),
 ('NoSQL', 0.3380617018914066),
 ('neural networks', 0.1889822365046136),
 ('deep learning', 0.1889822365046136),
 ('artificial intelligence', 0.1889822365046136),
 # ...
]
```

Выглядит довольно приличной рекомендацией для того, чьи интересы лежат в области Big Data и связаны с базами данных. (Весы не имеют какого-то особого смысла; они используются всего лишь для упорядочивания.)

Этот подход перестает справляться со своей задачей, когда число элементов становится очень большим. Напомним, что в связи с проблемой проклятия размерности, рассмотренной в *главе 12*, в многомерных векторных пространствах большинство векторов расположены друг от друга на большом расстоянии (и поэтому указывают в совершенно разных направлениях). Другими словами, при наличии большого количества интересующих тем "наиболее похожие пользователи" могут оказаться совсем непохожими на конкретного пользователя.

Возьмем сайт наподобие Amazon.com, где за последние пару десятилетий я сделал тысячи покупок. Можно попытаться установить похожих на меня пользователей на основе присущей мне модели покупательского поведения, но, скорее всего, во всем мире нет никого, чья покупательская история выглядит даже отдаленно похожей на мою. Кто бы ни был этот "наиболее похожий" на меня покупатель, он, вероятно, совсем не похож на меня, и, как результат, его покупки почти наверняка будут содействовать отвратительным рекомендациям.

## Коллаборативная фильтрация по схожести предметов

Альтернативный подход заключается в вычислении сходства непосредственно между интересующими темами. После этого можно сгенерировать рекомендации по каждому пользователю путем агрегирования тем, похожих на те, которые интересуют его в данное время.

Для начала следует *транспонировать* матрицу сфер интересов пользователей, благодаря чему строки соответствуют темам, а столбцы — пользователям:

```
# матрица групп пользователей по интересам в формате (тема, пользователи)
interest_user_matrix = [[user_interest_vector[j]
                        for user_interest_vector in user_interest_matrix]
                       for j, _ in enumerate(unique_interests)]
```

Как выглядит матрица групп пользователей по интересам `interest_user_matrix`? Строка `j` матрицы — это столбец `j` матрицы сфер интересов пользователей `user_interest_matrix`. Иными словами, она содержит 1 для пользователя, у которого есть интерес к данной теме, и 0 для пользователя, у которого его нет.

Например, `unique_interests[0]` равен `Big Data`, значит, `interest_user_matrix[0]` равен:

```
[1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0]
```

потому что пользователи 0, 8 и 9 указали свою заинтересованность в `Big Data`.

Теперь можно снова воспользоваться косинусным коэффициентом подобия. Если ровно те же самые пользователи заинтересованы в двух темах, то их сходство будет равно 1. Если нет двух пользователей, которых интересуют обе темы, то их коэффициент подобия будет равен 0.

```
# сходство интересов
```

```
interest_similarities = [[cosine_similarity(user_vector_i, user_vector_j)
                        for user_vector_j in interest_user_matrix]
                       for user_vector_i in interest_user_matrix]
```

Например, при помощи следующей функции можно найти темы, максимально похожие на тему `Big Data` (тема 0):

```
# темы, максимально похожие на тему interest_id
def most_similar_interests_to(interest_id):
    similarities = interest_similarities[interest_id]
    pairs = [(unique_interests[other_interest_id], similarity)
             for other_interest_id, similarity in enumerate(similarities)
             if interest_id != other_interest_id and similarity > 0]
    return sorted(pairs,
                  key=lambda (_, similarity): similarity,
                  reverse=True)
```

Функция сгенерирует следующие похожие темы:

```
[('Hadoop', 0.8164965809277261),
 ('Java', 0.6666666666666666),
 ('MapReduce', 0.5773502691896258),
 ('Spark', 0.5773502691896258),
 ('Storm', 0.5773502691896258),
 ('Cassandra', 0.4082482904638631),
 ('artificial intelligence', 0.4082482904638631),
 ('deep learning', 0.4082482904638631),
 ('neural networks', 0.4082482904638631),
 ('HBase', 0.3333333333333333)]
```

Теперь можно сгенерировать для пользователя рекомендации, просуммировав коэффициенты подобию тем, которые похожи на интересующие его темы:

```
# рекомендации на основе предметов
def item_based_suggestions(user_id, include_current_interests=False):
    # список рекомендаций, где суммируются похожие темы
    suggestions = defaultdict(float)
    # вектор тем, интересующих пользователя user_id
    user_interest_vector = user_interest_matrix[user_id]
    for interest_id, is_interested in enumerate(user_interest_vector):
        if is_interested == 1:
            # темы, похожие на тему interest_id
            similar_interests = most_similar_interests_to(interest_id)
            for interest, similarity in similar_interests:
                suggestions[interest] += similarity

    # упорядочить по весу
    suggestions = sorted(suggestions.items(),
                        key=lambda (_, similarity): similarity,
                        reverse=True)

    if include_current_interests:
        return suggestions
    else:
        return [(suggestion, weight)
                for suggestion, weight in suggestions
                if suggestion not in users_interests[user_id]]
```

Для пользователя 0 эта функция сгенерирует следующие (на вид разумные) рекомендации:

```
[('MapReduce', 1.861807319565799),
 ('Postgres', 1.3164965809277263),
 ('MongoDB', 1.3164965809277263),
 ('NoSQL', 1.2844570503761732),
 ('programming languages', 0.5773502691896258),
 ('MySQL', 0.5773502691896258),
 ('Haskell', 0.5773502691896258),
 ('databases', 0.5773502691896258),
 ('neural networks', 0.4082482904638631),
 ('deep learning', 0.4082482904638631),
 ('C++', 0.4082482904638631),
 ('artificial intelligence', 0.4082482904638631),
 ('Python', 0.2886751345948129),
 ('R', 0.2886751345948129)]
```

## Для дальнейшего изучения

- ◆ **Crab** (<http://muricoca.github.io/crab>) — это прикладная среда для создания рекомендательных систем (recommender systems) на языке Python.
- ◆ Библиотека **Graphlab** также располагает инструментарием для создания рекомендательных систем (<https://turi.com/products/create/docs/graphlab.toolkits.recommender.html>).
- ◆ Премия **Netflix** (<http://www.netflixprize.com>) — это достаточно известный конкурс по созданию усовершенствованной системы предложения рекомендаций о фильмах пользователям сети потокового мультимедиа Netflix.

# Базы данных и SQL

Память — лучший друг и худший враг человека.

Гилберт Паркер<sup>1</sup>

Необходимые данные часто размещены в системах управления базами данных (СУБД), предназначенных для эффективного хранения и извлечения данных. Основная их масса является реляционными, такими как Oracle, MySQL и SQL Server, в которых данные хранятся в виде таблиц, а запросы к которым, как правило, выполняются при помощи декларативного языка управления данными — языка структурированных запросов (Structured Query Language, SQL).

Язык SQL является существенной частью инструментария аналитика данных. В этой главе будет разработана питоновская реализация "не совсем базы данных" — приложение NotQuiteABase. Будут рассмотрены основы запросов к базам данных на языке SQL на примере этого приложения, которое представляет собой демонстрацию с самого что ни на есть абсолютно чистого листа, которую можно было придумать, принципа работы языка запросов. Надеюсь, что работа с NotQuiteABase даст хорошее представление о том, как решаются аналогичные задачи при помощи SQL.

## Операторы *CREATE TABLE* и *INSERT*

Реляционная база данных — это набор таблиц (и связей между ними). Таблица — это просто набор строк, похожий на матрицу, с которыми мы работали ранее. Однако таблица имеет связанную с ней фиксированную *схему*, состоящую из имен и типов столбцов.

Например, пусть дан набор данных `users`, содержащий идентификатор `user_id`, имя `name` и число друзей `num_friends` для каждого пользователя:

```
users = [[0, "Hero", 0],
         [1, "Dunn", 2],
         [2, "Sue", 3],
         [3, "Chi", 3]]
```

На языке SQL такую таблицу можно создать при помощи следующего оператора:

```
CREATE TABLE users (
    user_id INT NOT NULL,
```

---

<sup>1</sup> Гилберт Паркер (1862–1932) — канадский новеллист и английский политик. — *Прим. пер.*

```
name VARCHAR(200),
num_friends INT);
```

Следует обратить внимание, в операторе указано, что поля `user_id` и `num_friends` должны быть целыми числами (и что `user_id` не может быть `NULL`, обозначая отсутствующее значение наподобие `None` в языке Python) и что имя должно быть строкой длиной не более 200 символов. Приложение `NotQuiteABase` не учитывает типы, но мы представим, что оно это делает.



Язык SQL практически полностью независим от регистра букв и отступов. Прописные буквы и отступы в данном случае являются моим предпочтительным стилем. Если вы начинаете изучение SQL, то наверняка встретите примеры, которые стилизованы по-другому.

Вставить строки можно при помощи оператора `INSERT`:

```
INSERT INTO users (user_id, name, num_friends) VALUES (0, 'Hero', 0);
```

Кроме того, следует обратить внимание, что операторы SQL должны заканчиваться точкой с запятой, и SQL требует одинарные кавычки для строк.

В приложении `NotQuiteABase` экземпляр класса `Table`, который представляет таблицу данных, создается путем указания имен ее столбцов. И чтобы вставить строку, используется метод `insert()` класса `Table`, принимающий список значений отдельной строки таблицы, которые должны быть в том же порядке, что и имена столбцов.

Структурно каждая строка будет храниться как словарь из имен столбцов и их значений. В настоящей СУБД такое неэффективное по занимаемой памяти представление не используется, однако оно намного упрощает работу с приложением `NotQuiteABase`:

```
class Table:
    def __init__(self, columns):
        self.columns = columns
        self.rows = []

    def __repr__(self):
        """приемлемое представление таблицы: столбцы, затем строки"""
        return str(self.columns) + "\n" + "\n".join(map(str, self.rows))

    def insert(self, row_values):
        if len(row_values) != len(self.columns):
            raise TypeError("неправильное число элементов")
        row_dict = dict(zip(self.columns, row_values))
        self.rows.append(row_dict)
```

Например, можно задать:

```
users = Table(["user_id", "name", "num_friends"])
users.insert([0, "Hero", 0])
users.insert([1, "Dunn", 2])
```

```
users.insert([2, "Sue", 3])
users.insert([3, "Chi", 3])
users.insert([4, "Thor", 3])
users.insert([5, "Clive", 2])
users.insert([6, "Hicks", 3])
users.insert([7, "Devin", 2])
users.insert([8, "Kate", 2])
users.insert([9, "Klein", 3])
users.insert([10, "Jen", 1])
```

Если теперь напечатать `print(users)`, то получим:

```
['user_id', 'name', 'num_friends']
{'user_id': 0, 'name': 'Hero', 'num_friends': 0}
{'user_id': 1, 'name': 'Dunn', 'num_friends': 2}
{'user_id': 2, 'name': 'Sue', 'num_friends': 3}
...
```

## Оператор **UPDATE**

Иногда необходимо обновить данные, которые уже есть в базе данных. Например, если Dunn приобретает еще одного друга, то, чтобы отразить этот факт, понадобится такой код:

```
UPDATE users
SET num_friends = 3
WHERE user_id = 1;
```

Ключевые свойства обновления следующие:

- ◆ какую таблицу обновлять;
- ◆ какие строки обновлять;
- ◆ какие поля обновлять;
- ◆ какими должны быть их новые значения.

Добавим в `NotQuiteABase` аналогичный метод `update()`. Его первым аргументом будет словарь `dict`, ключами которого являются столбцы для обновления, а значениями — новые значения для этих полей. Вторым аргументом будет предикативная функция `predicate`, которая возвращает `True` для строк, которые должны быть обновлены, и в противном случае — `False`:

```
def update(self, updates, predicate):
    for row in self.rows:
        if predicate(row):
            for column, new_value in updates.iteritems():
                row[column] = new_value
```

Теперь можно сделать так:

```
users.update({'num_friends': 3}, # установить num_friends=3
             lambda row: row['user_id'] == 1) # в строках, где user_id==1
```

## Оператор *DELETE*

Есть два способа удаления строк из таблицы в SQL. Опасный способ удаляет все строки таблицы:

```
DELETE FROM users;
```

При менее опасном способе добавляется условный оператор `WHERE` и удаляются только те строки, которые удовлетворяют некоторому условию:

```
DELETE FROM users WHERE user_id = 1;
```

Добавить этот функционал к классу `Table` очень легко:

```
def delete(self, predicate=lambda row: True):
    """удалить все строки, удовлетворяющие предикату,
    или все строки, если предикативная функция не предоставлена"""
    self.rows = [row for row in self.rows if not(predicate(row))]
```

Если предоставить предикативную функцию (т. е. условный оператор `WHERE`), то будут удалены только те строки, которые удовлетворяют условию. При ее отсутствии по умолчанию предикативная функция всегда возвращает `True`, и будут удалены все строки.

Например:

```
users.delete(lambda row: row["user_id"] == 1) # удаляет строки
                                                # с user_id == 1
users.delete() # удаляет все строки
```

## Оператор *SELECT*

Обычно таблицы SQL не просматриваются непосредственно. Вместо этого они опрашиваются при помощи оператора `SELECT`:

```
SELECT * FROM users;           -- получить все содержимое
SELECT * FROM users LIMIT 2;  -- получить первые две строки
SELECT user_id FROM users;    -- взять только отдельные столбцы
SELECT user_id FROM users WHERE name = 'Dunn'; -- взять только
                                                -- отдельные строки
```

Оператор `SELECT` также можно использовать для вычисления полей:

```
SELECT LENGTH(name) AS name_length FROM users;
```

Предоставим классу `Table` метод `select()`, который возвращает новую таблицу. Метод принимает два необязательных аргумента:

- ◆ `keep_columns` задает названия столбцов, которые требуется сохранить в результирующей таблице. Если он не указан, то результат будет содержать все столбцы;
- ◆ `additional_columns` — это словарь, ключами которого являются новые имена столбцов, а значения — функциями, определяющими порядок вычисления значений новых столбцов.



Если не предоставить ни одного из них, то будет получена копия таблицы:

```
def select(self, keep_columns=None, additional_columns=None):

    if keep_columns is None:          # если столбцы не указаны,
        keep_columns = self.columns  # вернуть все столбцы

    if additional_columns is None:
        additional_columns = {}

    # новая таблица для результатов
    result_table = Table(keep_columns + additional_columns.keys())

    for row in self.rows:
        new_row = [row[column] for column in keep_columns]
        for column_name, calculation in additional_columns.iteritems():
            new_row.append(calculation(row))
        result_table.insert(new_row)

    return result_table
```

Метод `select()` возвращает новую таблицу, в то время как обычная операция SQL `SELECT` всего лишь генерирует некий промежуточный результирующий набор (если явно не вставлять результаты в таблицу).

Помимо этого, потребуются методы `where()` и `limit()`. Реализация обеих функций достаточно проста:

```
def where(self, predicate=lambda row: True):
    """вернуть только те строки, которые удовлетворяют
    указанному предикату"""
    where_table = Table(self.columns)
    where_table.rows = filter(predicate, self.rows)
    return where_table
```

```
def limit(self, num_rows):
    """вернуть только первые num_rows строк"""
    limit_table = Table(self.columns)
    limit_table.rows = self.rows[:num_rows]
    return limit_table
```

Теперь можно легко сконструировать функции `NotQuiteABase`, эквивалентные операторам SQL, указанным в комментариях перед ними:

```
# SELECT * FROM users;
users.select()

# SELECT * FROM users LIMIT 2;
users.limit(2)
```

```
# SELECT user_id FROM users;
users.select(keep_columns=["user_id"])

# SELECT user_id FROM users WHERE name = 'Dunn';
users.where(lambda row: row["name"] == "Dunn") \
    .select(keep_columns=["user_id"])

# SELECT LENGTH(name) AS name_length FROM users;
def name_length(row): return len(row["name"])

users.select(keep_columns=[],
             additional_columns = { "name_length" : name_length })
```

Обратим внимание — в отличие от остальной части книги — здесь использована обратная косая черта "\" для продолжения операторов на нескольких строках. По моему мнению, такое написание делает связанные в цепочку запросы NotQuiteABase читаемее любого другого способа их написания.

## Оператор GROUP BY

Другим распространенным оператором SQL является GROUP BY, который группирует строки с одинаковыми значениями в указанных полях и генерирует агрегатные значения, такие как MIN, MAX, COUNT и SUM. (Напоминают функцию группировки по полю group\_by из разд. "Управление данными" главы 10.)

Например, может понадобиться найти число пользователей и наименьший идентификатор пользователя user\_id для каждой возможной длины имени:

```
SELECT LENGTH(name) AS name_length,
       MIN(user_id) AS min_user_id,
       COUNT(*)     AS num_users
FROM users
GROUP BY LENGTH(name);
```

Каждое выбираемое оператором SELECT поле должно быть сгруппировано либо на основе условного оператора GROUP BY (каким является name\_length), либо на основе агрегатной функции (которыми являются min\_user\_id и num\_users).

Кроме того, язык SQL поддерживает условный оператор HAVING, который ведет себя аналогично условному оператору WHERE, за исключением того, что его фильтр применяется к агрегатным функциям (в отличие от него выражение WHERE фильтрует строки до того, как состоялось агрегирование).

Может понадобиться узнать среднее число друзей для пользователей, чьи имена начинаются с определенных букв, но увидеть только результаты для тех букв, чье соответствующее среднее больше 1. (Надо признать, что некоторые приводимые примеры надуманны.)

```
SELECT SUBSTR(name, 1, 1) AS first_letter,
       AVG(num_friends)  AS avg_num_friends
```

```
FROM users
GROUP BY SUBSTR(name, 1, 1)
HAVING AVG(num_friends) > 1;
```

(Функции для работы со строками разнятся от реализации к реализации SQL; в место этого в некоторых СУБД может использоваться функция SUBSTRING либо что-то в этом роде.)

Помимо этого, можно вычислять суммарные количества. В этом случае выражение GROUP BY опускается:

```
SELECT SUM(user_id) AS user_id_sum
FROM users
WHERE user_id > 1;
```

Чтобы реализовать такой функционал в класс Table приложения NotQuiteABase, добавим метод group\_by(). Он принимает имена столбцов, которые требуется сгруппировать, словарь агрегирующих функций, которые требуется применить к каждой группе, и необязательный предикат having, который обрабатывает несколько строк.

Затем он выполняет следующие шаги:

1. Создает словарь defaultdict для отображения кортежей tuple (со значениями group-by) в строки (содержащие значения group-by). Напомним, что использовать списки в качестве ключей словаря dict нельзя; следует использовать кортежи.
2. Выполняет обход строк таблицы, заполняя defaultdict.
3. Создает новую таблицу с правильными столбцами на выходе.
4. Выполняет обход словаря defaultdict и заполняет результирующую таблицу, применяя фильтр having, если он задан.

(Настоящая СУБД почти наверняка делает это более эффективным способом.)

```
# функция группировки
def group_by(self, group_by_columns, aggregates, having=None):
    grouped_rows = defaultdict(list)

    # заполнить группы
    for row in self.rows:
        key = tuple(row[column] for column in group_by_columns)
        grouped_rows[key].append(row)

    # результирующая таблица состоит из столбцов group_by и агрегатов
    result_table = Table(group_by_columns + aggregates.keys())

    for key, rows in grouped_rows.iteritems():
        if having is None or having(rows):
            new_row = list(key)
            for aggregate_name, aggregate_fn in aggregates.iteritems():
                new_row.append(aggregate_fn(rows))
            result_table.insert(new_row)

    return result_table
```

Снова посмотрим, как создать эквивалент предыдущих операторов SQL. Метрические показатели `name_length` следующие:

```
def min_user_id(rows): return min(row["user_id"] for row in rows)
```

```
stats_by_length = users \
    .select(additional_columns={"name_length" : name_length}) \
    .group_by(group_by_columns=["name_length"],
              aggregates={"min_user_id" : min_user_id,
                          "num_users" : len })
```

Метрические показатели для первой буквы `first_letter`:

```
def first_letter_of_name(row):
    return row["name"][0] if row["name"] else ""
```

```
def average_num_friends(rows):
    return sum(row["num_friends"] for row in rows) / len(rows)
```

```
def enough_friends(rows):
    return average_num_friends(rows) > 1
```

```
avg_friends_by_letter = users \
    .select(additional_columns={'first_letter' : first_letter_of_name}) \
    .group_by(group_by_columns=['first_letter'],
              aggregates={"avg_num_friends" : average_num_friends },
              having=enough_friends)
```

и сумма по идентификатору пользователя `user_id_sum`:

```
def sum_user_ids(rows): return sum(row["user_id"] for row in rows)
```

```
user_id_sum = users \
    .where(lambda row: row["user_id"] > 1) \
    .group_by(group_by_columns=[],
              aggregates={"user_id_sum" : sum_user_ids })
```

## Оператор **ORDER BY**

Часто требуется отсортировать результаты. Например, может понадобиться получить (в алфавитном порядке) первые два имени пользователей:

```
SELECT * FROM users
ORDER BY name
LIMIT 2;
```

Это легко реализовать, если в классе `Table` для таблиц есть метод `order_by()`, который принимает функцию `order`:

```
def order_by(self, order):
    new_table = self.select() # сделать копию
```

```
new_table.rows.sort(key=order)
return new_table
```

которую затем можно применить следующим образом:

```
friendliest_letters = avg_friends_by_letter \
    .order_by(lambda row: -row["avg_num_friends"]) \
    .limit(4)
```

Выражение SQL ORDER BY позволяет указывать порядок сортировки ASC (в возрастающем порядке) или DESC (в убывающем порядке) для каждого сортируемого поля; все это должно разместиться в функции order.

## Оператор JOIN

Таблицы реляционных баз данных часто *нормализуют*, т. е. их организуют так, чтобы уменьшить избыточность. Например, работая на языке Python с интересующими пользователей темами, можно предоставить каждому пользователю список только с его темами.

Таблицы SQL, как правило, не могут содержать списки, поэтому типичным решением является создание второй таблицы user\_interests, содержащей отношение "один-ко-многим" между идентификаторами пользователей user\_id и интересующими темами. На SQL это будет так:

```
CREATE TABLE user_interests (
    user_id INT NOT NULL,
    interest VARCHAR(100) NOT NULL
);
```

Тогда как в приложении NotQuiteABase надо создать таблицу:

```
user_interests = Table(["user_id", "interest"])
user_interests.insert([0, "SQL"])
user_interests.insert([0, "NoSQL"])
user_interests.insert([2, "SQL"])
user_interests.insert([2, "MySQL"])
```



До сих пор остается еще много избыточности — интересующая тема "SQL" хранится в двух разных местах. В реальной базе данных можно хранить идентификаторы user\_id и interest\_id в таблице user\_interests, а затем создать третью таблицу interests, связав поля interest\_id с interest, чтобы каждое название интересующей темы хранилось всего один раз. Здесь, это неоправданно усложнило бы примеры.

Как анализировать данные, если они разбросаны по разным таблицам? Это делается путем объединения таблиц на основе оператора JOIN. Он сочетает строки в левой таблице с соответствующими строками в правой таблице, где значение слова "соответствующий" зависит от того, как задается объединение.

Например, чтобы отыскать пользователей, заинтересованных в SQL, делается следующий запрос:

```

SELECT users.name
FROM users
JOIN user_interests
  ON users.user_id = user_interests.user_id
WHERE user_interests.interest = 'SQL'

```

Оператор `JOIN` означает, что для каждой строки в таблице `users` нужно обратиться к идентификатору `user_id` и ассоциировать эту строку с каждой строкой в таблице `user_interests`, содержащей тот же самый идентификатор `user_id`.

Следует обратить внимание, что при этом указывается, какие таблицы объединять (`JOIN`) и по каким столбцам (`ON`). Это так называемый *оператор внутреннего объединения* (`INNER JOIN`), который возвращает сочетания (и только те сочетания) строк, которые соответствуют заданным критериям объединения.

Кроме того, есть оператор левого объединения `LEFT JOIN`, который в дополнение к сочетанию совпадающих строк возвращает строку для каждой строки левой таблицы с несовпадающими строками (в этом случае, все поля, которые приходят из правой таблицы, равны `NULL`).

Используя оператор `LEFT JOIN`, легко подсчитать количество интересующих тем для каждого пользователя:

```

SELECT users.id, COUNT(user_interests.interest) AS num_interests
FROM users
LEFT JOIN user_interests
  ON users.user_id = user_interests.user_id

```

Операция левого объединения `LEFT JOIN` гарантирует, что пользователи, не имеющие интересующих их тем, по-прежнему будут иметь строки в объединенном наборе данных (с пустыми значениями полей `NULL` для полей из `user_interests`), и функция `COUNT` сосчитает только значения, которые не равны `NULL`.

Реализация функции `join()` в `NotQuiteABase` будет более ограниченной — она просто объединит две таблицы по любым общим столбцам. Написать ее даже в таком виде — задача не из тривиальных:

```

# функция объединения таблиц
def join(self, other_table, left_join=False):

    join_on_columns = [c for c in self.columns          # столбцы
                       if c in other_table.columns]   # в обеих таблицах

    additional_columns = [c for c in other_table.columns # столбцы только
                           if c not in join_on_columns] # в правой таблице

    # все столбцы из левой таблицы + дополнительные additional_columns
    # из правой
    join_table = Table(self.columns + additional_columns)

```

```

for row in self.rows:
    def is_join(other_row):
        return all(other_row[c] == row[c] for c in join_on_columns)

    other_rows = other_table.where(is_join).rows

    # каждая строка, которая удовлетворяет предикату,
    # генерирует результирующую строку
    for other_row in other_rows:
        join_table.insert([row[c] for c in self.columns] +
                           [other_row[c] for c in additional_columns])

    # если ни одна строка не удовлетворяет условию и это
    # левое объединение left join, то сгенерировать строку из None
    if left_join and not other_rows:
        join_table.insert([row[c] for c in self.columns] +
                           [None for c in additional_columns])

return join_table

```

Теперь можно найти пользователей, которые заинтересованы в SQL:

```

sql_users = users \
    .join(user_interests) \
    .where(lambda row: row["interest"] == "SQL") \
    .select(keep_columns=["name"])

```

И можно подсчитать интересующие темы:

```

def count_interests(rows):
    """Подсчитывает, сколько строк имеют ненулевые интересующие темы"""
    return len([row for row in rows if row["interest"] is not None])

user_interest_counts = users \
    .join(user_interests, left_join=True) \
    .group_by(group_by_columns=["user_id"],
              aggregates={"num_interests": count_interests })

```

В SQL есть еще ряд аналогичных операций: операция *правого объединения* RIGHT JOIN хранит строки из правой таблицы, для которых нет совпадений; операция *полного внешнего объединения* FULL OUTER JOIN хранит строки из обеих таблиц, для которых нет совпадений. Они останутся без реализации.

## Подзапросы

В SQL результаты запросов можно обрабатывать при помощи операторов SELECT и JOIN точно так же, как если бы они были таблицами. Так, если бы нужно было найти наименьший идентификатор пользователя user\_id из тех, кто заинтересован

в SQL, то можно использовать *подзапрос*. (Естественно, можно сделать то же самое вычисление с помощью операции JOIN, но тогда не будут проиллюстрированы подзапросы.)

```
SELECT MIN(user_id) AS min_user_id FROM
(SELECT user_id FROM user_interests WHERE interest = 'SQL') sql_interests;
```

Учитывая то, как разработано приложение NotQuiteABase, мы получаем это бесплатно. (Результаты запросов фактически представлены таблицами.)

```
likes_sql_user_ids = user_interests \
    .where(lambda row: row["interest"] == "SQL") \
    .select(keep_columns=['user_id'])

likes_sql_user_ids.group_by(group_by_columns=[],
                             aggregates={"min_user_id" : min_user_id })
```

## Индексы

Чтобы найти строки, содержащие определенное значение (допустим, с именем "Hero"), приложение NotQuiteABase должно проверять каждую строку в таблице. Если таблица имеет много строк, то поиск может занять очень много времени.

То же касается и алгоритма объединения join — он крайне неэффективен. Для каждой строки в левой таблице в поисках совпадений он просматривает каждую строку в правой таблице. Для двух больших таблиц это может занять вечность.

Часто требуется применять ограничения к некоторым столбцам. Например, в таблице users может потребоваться, чтобы у двух различных пользователей не было одинаковых идентификаторов user\_id.

Все эти задачи решают *индексы*. Если бы в таблице user\_interests имелся индекс по полю user\_id, то умный алгоритм объединения join нашел бы совпадения непосредственно, а не в результате просмотра всей таблицы. Если бы в таблице users был "уникальный" индекс по полю user\_id, то при попытке вставить дубликат было бы выведено сообщение об ошибке.

Каждая таблица в базе данных может иметь один или несколько индексов, которые позволяют быстро выполнять просмотр таблиц по ключевым столбцам, эффективно объединять таблицы и накладывать уникальные ограничения на столбцы или их сочетания.

Конструирование и применение индексов часто относят скорее к приемам (количество которых разнится в зависимости от конкретной СУБД) квалификации высокого уровня, и если работа будет связана с выполнением большого числа запросов к базе данных, то с этими приемами следует познакомиться.



## Оптимизация запросов

Напомним запрос для поиска всех пользователей, которые интересуются SQL:

```
SELECT users.name
FROM users
JOIN user_interests
  ON users.user_id = user_interests.user_id
WHERE user_interests.interest = 'SQL'
```

В приложении NotQuiteABase можно написать такой запрос (как минимум) двумя различными способами. Можно отфильтровать таблицу `user_interests` перед тем, как выполнить объединение:

```
user_interests \
  .where(lambda row: row["interest"] == "SQL") \
  .join(users) \
  .select(["name"])
```

Либо отфильтровать результаты объединения:

```
user_interests \
  .join(users) \
  .where(lambda row: row["interest"] == "SQL") \
  .select(["name"])
```

В обоих случаях в конечном счете результаты будут одинаковыми, но фильтрация до объединения более эффективна, поскольку в этом случае для обработки функцией `join` остается намного меньше строк.

В SQL можно вообще не беспокоиться об этом. Нужно только "объявить" результаты, которые требуется получить, и передать их ядру выполнения запросов для обработки (и эффективного использования индексов).

## Базы данных NoSQL

В последние годы наметилась тенденция применения нереляционных (слабоструктурированных) СУБД "NoSQL", в которых данные представлены не таблицами. Например, MongoDB — это популярная СУБД, в которой схемы данных не используются, а ее элементы не строки, а документы в формате JSON произвольной сложности.

Существуют самые разнообразные виды нереляционных баз данных: базы данных, в которых данные хранятся не в строках, а в столбцах (они подходят для тех случаев, когда есть много столбцов, а запросы выполняются только к нескольким из них); хранилища данных в формате "ключ-значение", которые оптимальны в тех случаях, когда нужно получить единое (комплексное) значение по ключу; графовые базы данных для хранения и поиска в графах; базы данных, оптимизированные для работы между несколькими центрами обработки данных; базы данных, предназна-

ченые для работы в оперативной памяти; базы данных для хранения данных временных рядов, и сотни других.

Характерный тренд на ближайшее будущее, возможно, еще не оформился, поэтому пока ограничимся лишь сообщением о существовании СУБД NoSQL, чтобы быть в курсе о таком способе хранения данных.

## Для дальнейшего изучения

- ◆ Тем, кому нужна реляционная СУБД, сообщая несколько вариантов, которые доступны для скачивания: быстрая и миниатюрная СУБД SQLite (<http://www.sqlite.org/>), более крупные и с большим набором характеристик СУБД MySQL (<http://www.mysql.com/>) и PostgreSQL (<https://www.postgresql.org/>). Все они бесплатные и располагают обширной документацией.
- ◆ Тем, кто желает заняться исследованием СУБД NoSQL, лучше всего начать с простой MongoDB (<https://www.mongodb.com/>), которая может стать как благословением, так и своего рода проклятием. Кроме того, она располагает достаточно хорошей документацией.
- ◆ Статья в Википедии, посвященная СУБД NoSQL (<https://en.wikipedia.org/wiki/NoSQL>), почти наверняка сейчас содержит ссылки на системы управления базами данных, которые даже не существовали на момент написания этой книги.

# Распределенные вычисления MapReduce

Будущее уже наступило. Просто оно еще не распределено равномерно.

Уильям Гибсон<sup>1</sup>

MapReduce — это модель программирования для выполнения параллельной обработки больших наборов данных. Несмотря на то, что она представляет собой мощную технологию, ее основания достаточно просты.

Представим, что имеется совокупность элементов, которые как-то требуется обработать. Например, элементами могут быть журналы сайта, тексты различных книг, файлы изображений или что-нибудь еще. Базовая версия алгоритма MapReduce состоит из следующих шагов:

1. Применить функцию трансформации `mapper` для преобразования каждого элемента в ноль или более пар в формате "ключ-значение". (Если быть точным, эта функция называется `map`, однако в Python уже есть функция с таким названием, и чтобы не смешивать два понятия, возьмем за основу первое.)
2. Собрать вместе все пары с одинаковыми ключами.
3. Использовать функцию свертки `reducer` к каждой совокупности сгруппированных значений для создания результирующих значений для соответствующего ключа.

Все это выглядит абстрактно, поэтому обратимся к конкретному примеру. В науке о данных существует несколько абсолютных правил, и одно из них заключается в том, что первый пример работы алгоритма MapReduce должен быть связан с подсчетом частотности слов.

## Пример: подсчет частотности слов

Социальная сеть DataSciencester выросла до миллионов пользователей! Великолепная гарантия обеспеченности работой, правда, это несколько осложняет выполнение рутинных обязанностей по анализу данных.

К примеру, директор по контенту хотел бы знать, о чем люди пишут в своих постах, размещаемых в ленте новостей (`status update`). В качестве первой попытки вы

---

<sup>1</sup> Уильям Форд Гибсон (1948) — американский писатель-фантаст. Считается основателем стиля киберпанк, определившего жанровое лицо литературы 1980-х гг. (см. [https://ru.wikipedia.org/wiki/Гибсон,\\_Уильям](https://ru.wikipedia.org/wiki/Гибсон,_Уильям)). — *Прим. пер.*

решаете подсчитать слова, которые встречаются в постах, чтобы подготовить отчет о наиболее часто употребляемых словах.

При наличии нескольких сотен пользователей это сделать несложно:

```
# подсчет частотности слов (старая версия)
def word_count_old(documents):
    """подсчет частотности слов без использования технологии MapReduce"""
    return Counter(word
                    for document in documents
                    for word in tokenize(document))
```

Но когда пользователей — миллионы, совокупность документов `documents` (посты в ленте новостей) неожиданно становится слишком большой, чтобы поместиться на одном компьютере. Если бы удалось как-то разместить их в модели `MapReduce`, то можно было бы воспользоваться элементами инфраструктуры "больших данных", которую реализовали ваши инженеры.

Во-первых, потребуется функция, которая превращает документ в последовательность пар "ключ-значение" (назовем эту функцию проектором). При этом выходящие значения должны быть сгруппированы по слову, и, следовательно, ключи будут словами. И при появлении каждого слова будет генерироваться 1, обозначая, что эта пара соответствует одному вхождению слова в текст:

```
# проектор или функция трансформации на этапе Map
def wc_mapper(document):
    """для каждого слова в документе, генерировать (слово,1)"""
    for word in tokenize(document):
        yield (word, 1)
```

На время пропуска "механику" реализации 2-го этапа, представим, что для некоторого слова собран список из единиц. Тогда для получения суммы чисел для этого слова всего лишь нужно их просуммировать в функции свертки (назовем эту функцию редуктором):

```
# редуктор или функция свертки на этапе Reduce
def wc_reducer(word, counts):
    """суммировать числа для слова"""
    yield (word, sum(counts))
```

Вернемся ко 2-му этапу. Теперь необходимо собрать результаты работы проектора `wc_mapper` и передать их в редуктор `wc_reducer`. Подумаем, как это можно сделать, используя лишь один компьютер:

```
# подсчет частот слов
def word_count(documents):
    """подсчитать слова во входящих документах при помощи MapReduce"""

    # место для хранения сгруппированных значений
    collector = defaultdict(list)
```

```

for document in documents:
    for word, count in wc_mapper(document):
        collector[word].append(count)

return [output
        for word, counts in collector.iteritems()
        for output in wc_reducer(word, counts)]

```

Пусть имеется три документа ["data science", "big data", "science fiction"].

Тогда проектор `wc_mapper`, примененный к первому документу, вернет две пары ("data", 1) и ("science", 1). После проверки всех трех документов, переменная `collector` будет содержать следующее:

```

{ "data" : [1, 1],
  "science" : [1, 1],
  "big" : [1],
  "fiction" : [1] }

```

Затем редуктор `wc_reducer` сгенерирует итоговую частоту каждого слова:

```
[("data", 2), ("science", 2), ("big", 1), ("fiction", 1)]
```

## Почему MapReduce?

Как уже упоминалось ранее, основное преимущество технологии MapReduce заключается в том, что она позволяет распределять вычисления путем перемещения обработки в сторону данных. Представим, что приложение `wordcount` подсчитывает частотности слов в миллиардах документов.

Первоначальный подход (без применения технологии MapReduce) требует, чтобы компьютер выполнял обработку, имея доступ к любому документу. Следовательно, все документы должны либо располагаться на том же компьютере, либо передаваться ему во время обработки. И при этом, что представляется более важным, компьютер способен обрабатывать только по одному документу за один момент времени.



Возможно даже, что он сможет обрабатывать до нескольких документов одновременно, если у него несколько ядер и если переписать код, чтобы воспользоваться ими. Но даже в этом случае все документы все равно должны быть доставлены к этому компьютеру.

Теперь представим, что эти миллиарды документов распределены между 100 компьютеров. При наличии правильной инфраструктуры (и умолчав о некоторых подробностях) можно сделать следующее:

- ◆ на каждом компьютере прогнать документы через проектор с целью получения множества пар в формате "ключ-значение";
- ◆ распределить эти пары среди определенного числа редукторов или "свертывающих" компьютеров, следя за тем, чтобы все пары, соответствующие любому конкретному ключу, оставались на одном компьютере;

- ◆ на каждом редукторе выполнить группировку пар по ключу и затем прогнать через редуктор каждое подмножество значений;
- ◆ вернуть пары в формате "ключ — выходные данные".

Самое удивительное заключается в том, что этот метод масштабируется по горизонтали. Если удвоить количество компьютеров, то (игнорируя определенные фиксированные затраты на работу системы MapReduce) вычисления должны выполняться примерно в два раза быстрее. Каждый компьютер с проектором должен будет выполнить только половину работы, и (предполагая, что уникальных ключей для дальнейшего распределения работы среди редукторов имеется в достаточном количестве) то же самое относится к компьютерам с редуктором.

## MapReduce в более общей реализации

Если подумать, то весь зависимый код для выполнения подсчета количества вхождений уникальных слов (частотностей слов) в предыдущем примере содержится в функциях `wc_mapper` и `wc_reducer`. Поэтому пара изменений даст намного более общую реализацию (которая по-прежнему будет выполняться на одном компьютере):

```
def map_reduce(inputs, mapper, reducer):
    """применяет MapReduce ко входящим данным, используя
    проектор mapper и редуктор reducer"""
    collector = defaultdict(list)

    for input in inputs:
        for key, value in mapper(input):
            collector[key].append(value)

    return [output
            for key, values in collector.iteritems()
            for output in reducer(key, values)]
```

И затем можно подсчитать частоты слов так:

```
word_counts = map_reduce(documents, wc_mapper, wc_reducer)
```

Такая реализация предоставляет гибкость при решении широкого спектра задач.

Прежде чем продолжить, заметим, что функция свертки `wc_reducer` всего лишь суммирует значения, соответствующие каждому ключу. От этого распространенного вида агрегации можно абстрагироваться:

```
# свернуть значения, используя функцию агрегации
def reduce_values_using(aggregation_fn, key, values):
    """сворачивает пару 'ключ-значение', применяя
    агрегирующую функцию aggregation_fn к значениям"""
    yield (key, aggregation_fn(values))
```

```
# функция свертки значений
```

```
def values_reducer(aggregation_fn):
```

```
"""превращает функцию (значения -> выход) в функцию свертки reducer,
которая преобразует (ключ, значение) -> (ключ, выход)"""
return partial(reduce_values_using, aggregation_fn)
```

Теперь легко можно создать редуктор с агрегирующими функциями:

```
sum_reducer = values_reducer(sum)
max_reducer = values_reducer(max)
min_reducer = values_reducer(min)
count_distinct_reducer = values_reducer(lambda values: len(set(values)))
```

и т. д.

## Пример: анализ обновлений ленты новостей

Директора по контенту впечатлил отчет о частотности употреблений слов, и он интересуется, что еще можно узнать из обновлений ленты новостей? Вам удалось извлечь набор обновлений, который выглядят примерно так:

```
{"id": 1,
 "username" : "joelgrus",
 "text" : "Is anyone interested in a data science book?",
 "created_at" : datetime.datetime(2013, 12, 21, 11, 47, 0),
 "liked_by" : ["data_guy", "data_gal", "mike"] }
```

Допустим, необходимо выяснить, в какой день недели пользователи говорят о науке о данных чаще всего. Для того чтобы узнать это, просто подсчитаем, сколько в каждый день недели делалось обновлений, посвященных науке о данных. Иными словами, потребуется сгруппировать данные по дню недели, который и будет ключом. И если для каждого обновления, в котором содержится понятие "наука о данных", генерировать 1, то при помощи функции суммирования `sum` можно получить итоговое число упоминаний этого понятия:

```
# проектор для числа упоминаний термина науки о данных
def data_science_day_mapper(status_update):
    """для дня недели возвращает (day_of_week, 1), если
    обновление ленты status_update содержит "data science" """
    if "data science" in status_update["text"].lower():
        day_of_week = status_update["created_at"].weekday()
        yield (day_of_week, 1)
```

```
# дни, когда упоминалась наука о данных
data_science_days = map_reduce(status_updates,
                               data_science_day_mapper,
                               sum_reducer)
```

Возьмем более сложный пример и допустим, что по каждому пользователю нужно выяснить слово, которое он использует наиболее часто в своих обновлениях ленты новостей. На ум приходят три возможных подхода к реализации проектора `mapper`:





## Пример: умножение матриц

Вспомним из разд. "Умножение матриц" главы 21, что при наличии матрицы  $A$  размера  $m \times n$  и матрицы  $B$  размера  $n \times k$  можно их умножить и получить матрицу  $C$  размера  $m \times k$ , где элемент матрицы  $C$  в строке  $i$  и столбце  $j$  задается так:

$$c_{ij} = a_{i1}b_{1j} + a_{i2}b_{2j} + \dots + a_{in}b_{nj}.$$

Как было показано ранее, "естественным" способом представления матрицы размера  $m \times n$  является список списков, где элемент  $a_{ij}$  является  $j$ -м элементом  $i$ -го списка.

Однако большие матрицы иногда бывают разреженными, и вследствие этого, большинство их элементов равны нулю. Для больших разреженных матриц представление в виде списков списков может стать очень расточительным. Более компактным представлением является список кортежей (name, i, j, value), где имя name идентифицирует матрицу, а i, j, указывает на местоположение с ненулевым значением value.

Например, матрица размера миллиард на миллиард имеет *квинтиллион* записей, которые было бы непросто хранить на компьютере. Но если в каждой строке имеется лишь несколько ненулевых элементов, то данное альтернативное представление займет намного порядков меньше места.

При таком представлении можно воспользоваться технологией MapReduce для выполнения умножения матриц в распределенной среде.

В обоснование алгоритма следует обратить внимание на то, что каждый элемент  $a_{ij}$  используется только для вычисления элементов матрицы  $C$  в строке  $i$ , а каждый элемент  $b_{ij}$  — для вычисления элементов матрицы  $C$  в столбце  $j$ . Задача состоит в том, чтобы каждый результат редуктора reducer являлся единственной записью в матрице  $C$ , и, следовательно, понадобится, чтобы проектор mapper генерировал ключи, совпадающие с единственной записью в  $C$ . Это предполагает следующую реализацию:

```
# проектор для умножения матриц
def matrix_multiply_mapper(m, element):
    """m - это общая размерность (столбцы A, строки B)
       элемент - это кортеж (имя матрицы matrix_name, i, j,
                           значение value)"""
    name, i, j, value = element

    if name == "A":
        # A_ij - это j-я запись в сумме для каждого C_ik, k=1..m
        for k in range(m):
            # группировка с другими записями для C_ik
            yield((i, k), (j, value))
    else:
        # B_ij - это i-я запись в сумме для каждого C_kj
```

```

for k in range(m):
    # группировка с другими записями для C_kj
    yield((k, j), (i, value))

# редуктор для умножения матриц
def matrix_multiply_reducer(m, key, indexed_values):
    results_by_index = defaultdict(list)
    for index, value in indexed_values:
        results_by_index[index].append(value)

    # суммировать все произведения позиций с двумя результатами
    sum_product = sum(results[0] * results[1]
                       for results in results_by_index.values()
                       if len(results) == 2)

    if sum_product != 0.0:
        yield (key, sum_product)

```

Например, если имеются две матрицы:

```
A = [[3, 2, 0],
      [0, 0, 0]]
```

```
B = [[4, -1, 0],
      [10, 0, 0],
      [0, 0, 0]]
```

то их можно переписать в виде кортежей:

```

entries = [("A", 0, 0, 3), ("A", 0, 1, 2),
           ("B", 0, 0, 4), ("B", 0, 1, -1), ("B", 1, 0, 10)]
mapper = partial(matrix_multiply_mapper, 3)
reducer = partial(matrix_multiply_reducer, 3)

```

```
map_reduce(entries, mapper, reducer) # [(0, 1), -3], [(0, 0), 32]
```

Такая реализация не вызывает особого интереса на небольших матрицах, но при наличии миллионов строк и столбцов она будет целесообразной.

## Ремарка: сумматоры

Вдумчивый читатель, наверное, уже обратил внимание на то, что многие функции проекции mapper, по-видимому, содержат излишнюю информацию. Например, при подсчете частотностей слов вместо генерирования пар (word, 1) и суммирования их значений можно генерировать пары (word, None) и просто брать их длину.

Причина, почему этого не было сделано, заключается в том, что в распределенной среде иногда используются *сумматоры* (combiners), которые уменьшают объем данных, передаваемых от компьютера к компьютеру. Если один из компьютеров

с проектором видит слово "данные" 500 раз, то перед передачей результата компьютеру с редуктором он может объединить 500 пар ("data", 1) в одну пару ("data", 500). Это приводит к гораздо меньшему объему перемещаемых данных, что может сделать алгоритм существенно быстрее.

Реализованный выше редуктор позволяет обрабатывать эти объединенные данные правильно. (Если написать ее при помощи функции `len`, то этого не получится.)

## Для дальнейшего изучения

- ◆ Наиболее широко используемой системой на основе технологии MapReduce является Hadoop (<http://hadoop.apache.org/>), которая сама по себе достойна многих книг. С Hadoop связаны самые разные коммерческие и некоммерческие дистрибутивы и огромная экосистема инструментов.

Для того чтобы ею воспользоваться, необходимо создать собственный *кластер* (или найти поставщика, чтобы использовать предлагаемые им кластеры), что не обязательно является задачей для неслабонервных. Проекторы и редукторы Hadoop чаще всего написаны на Java, хотя существует механизм, известный как "поточная передача Hadoop", который позволяет писать их на других языках (в том числе на Python).

- ◆ Amazon.com предлагает веб-сервис Elastic MapReduce (<https://aws.amazon.com/ru/emr/>), который позволяет программно создавать и уничтожать кластеры, взимая плату только за количество времени, в течение которого они используются.
- ◆ `mrjob` (<https://github.com/Yelp/mrjob>) — это пакет на Python для взаимодействия с Hadoop (или с Elastic MapReduce от Amazon.com).
- ◆ Пакетные задания Hadoop, как правило, имеют высокую латентность и поэтому плохо подходят для анализа данных в реальном времени. Впрочем, есть различные инструменты для работы в "реальном времени", которые разработаны поверх Hadoop. Кроме того, появилось несколько альтернативных сред, которые становятся все более популярными. Двумя наиболее популярными из них являются Spark (<http://spark.apache.org/>) и Storm (<http://storm.apache.org/>).
- ◆ С учетом сказанного становится понятно, что в ближайшем будущем, вполне возможно, приобретут злободневность некоторые свежие распределенные системы обработки данных, которые еще не существовали на момент написания этой книги. Их можно отыскать самостоятельно.

# Идите и займитесь аналитикой

И сейчас, в очередной раз, я приглашаю мое ужасное потомство идти вперед и процветать.

*Мэри Шелли*<sup>1</sup>

Куда двигаться дальше? Если предположить, что эта книга не отпугнула читателя от науки о данных, то хочется отметить ряд тем, с которыми следует ознакомиться в следующую очередь.

## Интерактивная оболочка IPython

Интерактивная оболочка IPython (<http://ipython.org/>) уже ранее упоминалась в книге. Она обеспечивает гораздо больший функционал, чем стандартная среда Python, а также добавляет "волшебные" функции, которые позволяют (помимо всего прочего) легко копировать и вставлять код (что в обычных условиях осложняется сочетанием пустых строк и пробельных символов форматирования) и выполнять скрипты из оболочки.

Освоение IPython намного упростит работу. (Этого можно добиться, зная всего лишь несколько элементов интерфейса оболочки IPython.)

Кроме того, эта среда позволяет создавать "записные книжки", объединяющие текст, живой питоновский код и визуализации, которыми можно делиться с другими людьми или просто держать под рукой в виде дневника с записями о том, что было сделано (рис. 25.1)<sup>2</sup>:

---

<sup>1</sup> Мэри Шелли (1797–1851) — английская писательница, автор романа "Франкенштейн, или современный Прометей" (см. [https://ru.wikipedia.org/wiki/Шелли,\\_Мэри](https://ru.wikipedia.org/wiki/Шелли,_Мэри)). — *Прим. пер.*

<sup>2</sup> Записная книжка Python — это веб-приложение, предназначенное для интерактивного выполнения грамотных вычислений, в которых совмещены пояснительный текст, математические выкладки, вычисления и мультимедийное сопровождение. Входящие и выходящие данные хранятся в постоянных ячейках, которые могут редактироваться на месте. Простые текстовые документы, называемые записными книжками, служат для записи и распространения результатов богатых вычислений. Теперь IPython входит в состав приложения Jupyter Notebook App (<http://jupyter.org/>), объединяющего более 40 языков программирования. Примеры использования приложения: очистка и преобразование данных, численное и статистическое моделирование, машинное обучение и многое другое. — *Прим. пер.*

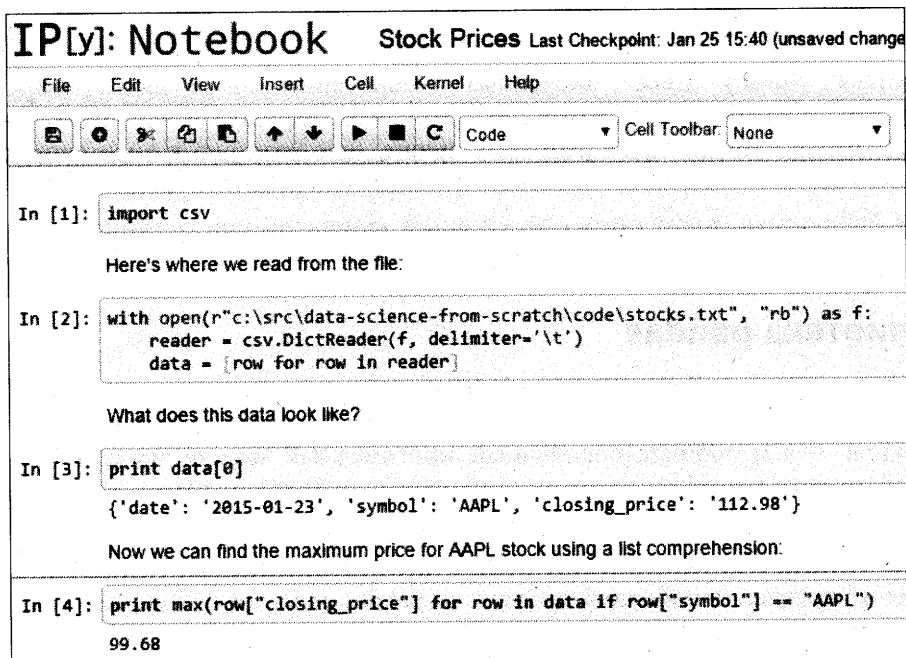


Рис. 25.1. Записная книжка IPython

## Математический аппарат

На протяжении всей книги мы поверхностно занимались линейной алгеброй (см. главу 4), математической статистикой (см. главу 5), теорией вероятностей (см. главу 6) и различными аспектами машинного обучения.

Чтобы стать хорошим специалистом в области анализа данных, требуется гораздо лучше разбираться в этих областях знаний, и я призываю вас заняться их более углубленным изучением, используя для этих целей учебники, рекомендованные в конце глав, учебники, которые вы предпочитаете сами, курсы дистанционного обучения и даже посещения реальных курсов.

## Не с чистого листа

Реализация методов "с чистого листа" позволяет лучше понять, как они работают. Но такой подход, как правило, сильно сказывается на их производительности (если только методы не реализуются специально для повышения производительности), простоте использования, быстроте прототипирования и обработке ошибок.

На практике лучше использовать хорошо продуманные библиотеки, прочно реализующие теоретические основы. (Мое первоначальное предложение для издательства, касающееся этой книги, подразумевало вторую часть "Изучение библиотек", на которую издательство O'Reilly, к счастью, наложило вето.)

## Библиотека NumPy

Библиотека NumPy (для "численных вычислений на Python", <http://www.numpy.org/>) обеспечивает условия для "реальных" научных вычислений. Она располагает массивами, которые работают лучше, чем матрицы в виде списков списков, и большим количеством числовых функций для работы с ними.

Данная библиотека может стать структурным элементом при создании других библиотек, что делает ее особенно ценной.

## Библиотека pandas

Библиотека pandas (<http://pandas.pydata.org/>) предоставляет дополнительные структуры данных для работы с массивами данных на языке Python. Ее основная абстракция — это проиндексированный многомерный массив значений DataFrame, который по сути похож на класс Table приложения NotQuiteABase, разработанного в главе 23, но с гораздо большим функционалом и лучшей производительностью.

Если вы собираетесь использовать Python для преобразования, разбиения, группирования и управления наборами данных, то pandas является бесценным инструментом для этих целей.

## Библиотека scikit-learn

Библиотека scikit-learn (<http://scikit-learn.org/>) — это, наверное, самая популярная библиотека для работы в области машинного обучения на языке Python. Она содержит все модели, которые были тут реализованы, и многие другие. В реальной ситуации не следует строить дерево принятия решений "с чистого листа"; всю тяжелую работу, связанную с решением этой задачи, должна делать библиотека scikit-learn. При решении реальных задач в области оптимизации вместо реализации какого-либо алгоритма оптимизации вручную следует положиться на библиотеку scikit-learn, где он уже эффективно реализован.

Документация по библиотеке содержит огромное количество примеров ([http://scikit-learn.org/stable/auto\\_examples/](http://scikit-learn.org/stable/auto_examples/)) того, что можно сделать с ее помощью (и того, какие вообще задачи решает машинное обучение).

## Визуализация

Диаграммы, построенные при помощи библиотеки matplotlib, получились ясными и функциональными, но не особо стилизованными (и совсем не интерактивными). Если вы хотите углубиться в область визуализации данных, то вот несколько вариантов.

В первую очередь стоит глубже исследовать библиотеку matplotlib. То, что было показано, является лишь незначительной частью ее функциональных возможностей. На сайте библиотеки содержится много примеров (<http://matplotlib.org/examples/>), которые их демонстрируют, а также галерея (<http://matplotlib.org/gallery.html>) из некоторых наиболее интересных. Если перед вами стоит задача

создания статических визуализаций (скажем, для книжных иллюстраций), то, вероятно, на следующем шаге лучше всего заняться углубленным изучением библиотеки `matplotlib`.

Помимо этого, стоит попробовать библиотеку `seaborn` (<http://web.stanford.edu/~mwaskom/software/seaborn/>), которая (среди прочего) делает диаграммы, `matplotlib` более привлекательными.

Если же стоит задача создания интерактивных визуальных презентаций, которыми можно делиться в сети, то очевидно лучше всего для этого подойдет библиотека `D3.js`, написанная на JavaScript и предназначенная для создания "документов, управляемых данными" (три буквы `D` от английского выражения `Data Driven Documents`). Даже если вы не знакомы с JavaScript, всегда можно взять примеры из галереи `D3` (<https://github.com/d3/d3/wiki/Gallery>) и настроить их для работы со своими данными. (Хорошие аналитики данных копируют примеры из галереи `D3`; великие аналитики их *заимствуют*.)

Даже если вы не заинтересованы в `D3`, то элементарный просмотр галереи сам по себе уже является невероятно поучительным занятием в области визуализации данных.

Проект `Bokeh` (<http://bokeh.pydata.org/>) способен привнести в Python функционал в стиле пакета `D3`.

## Язык программирования R

Хотя можно совершенно безнаказанно обойтись без овладения языком программирования `R` (<https://www.r-project.org/>), разработанного для статистической обработки данных и работы с графикой, во множестве проектов в области анализа данных и многими аналитиками он все же используется, поэтому стоит по крайней мере с ним ознакомиться.

Отчасти это нужно, чтобы появилась возможность понимать статьи, примеры и код в электронных журналах, ориентированных на пользователей языка `R`; отчасти, потому что знание языка поможет лучше оценить (сравнительно) чистую элегантность языка `Python`; а отчасти, поможет быть более информированным участником нескончаемых дискуссий по теме "`R` по сравнению с `Python`".

В мире нет недостатка в учебниках, курсах и книгах, посвященных языку `R`. Например, много и хорошо отзываюся о книге "Практическое программирование на `R`" ("`Hands-On Programming with R`", <http://shop.oreilly.com/product/0636920028574.do>), и не только потому, что это книга вышла в издательстве `O'Reilly`. (Ну, хорошо, в основном потому, что это книга вышла в издательстве `O'Reilly`.)

## Где найти данные?

Если решение задач в области науки о данных является частью вашей работы, то, скорее всего, вы получаете данные прямо на работе (впрочем, и не обязательно). Но как быть, если вы занимаетесь наукой о данных ради интеллектуального удовольствия? Хоть нас и окружают данные, вот несколько отправных точек.

- ◆ **Data.gov** (<https://www.data.gov/>) — это портал открытых данных правительства США. Если требуются данные по всему, что связано с государством (которые в наши дни приобретают особую значимость), то этот портал будет хорошей отправной точкой<sup>3</sup>.
- ◆ На социальном новостном сайте reddit есть пара форумов: `r/datasets`<sup>4</sup> и `r/data`<sup>5</sup>, где можно запросить и найти данные.
- ◆ **Amazon.com** поддерживает коллекцию общедоступных наборов данных (<https://aws.amazon.com/ru/public-data-sets/>), которые они хотели бы, чтобы пользователи применяли для анализа при помощи продуктов их производства (но, которые никто не мешает анализировать при помощи любых других продуктов по желанию).
- ◆ У Робба Ситона в его блоге (<http://rs.io/100-interesting-data-sets-for-statistics/>) есть причудливый список курируемых наборов данных.
- ◆ **Kaggle** (<https://www.kaggle.com/>) — это сайт, который проводит состязания в области науки о данных. Мне ни разу не удалось в них поучаствовать (в виду нехватки соревновательного духа, когда дело доходит до науки о данных), но вы могли бы попробовать.

## Занятия анализом данных

Конечно, замечательно, если умеешь просматривать каталоги с данными, и тем не менее лучшие проекты (и результаты) — это, наверное, те, которые были вызваны неким непреодолимым стремлением их реализовать. Вот несколько из тех, которые я выполнил.

## Новости хакера

Новости хакера (<https://news.ycombinator.com/news>) — это агрегатор новостей и дискуссионная площадка для новостей, связанных с технологиями. Сайт собирает огромное количество статей, многие из которых для меня не интересны.

Поэтому несколько лет назад я задался целью создать классификатор статей этого сайта (<https://github.com/joelgrus/hackernews>), который бы предсказывал, будет ли мне интересно то или иное повествование. С этим сайтом получилось не все так просто. Владельцы сайта отвергли саму идею о том, что кто-то может не быть заинтересован в любом из рассказов, выложенных на их сайте.

Работа включала в себя разметку в ручном режиме большого количества статей (с целью получения обучающей выборки), выбор характерных признаков статей

---

<sup>3</sup> В Российской Федерации тоже существует сайт правительства РФ с открытыми данными — <http://government.ru>. — *Ред.*

<sup>4</sup> <https://www.reddit.com/r/datasets>.

<sup>5</sup> <https://www.reddit.com/r/data>.



(например, слов в заголовке и доменных имен в ссылках), обучение наивного байесовского классификатора, аналогичного рассмотренному в книге спам-фильтру.

По причинам, ныне затерявшимся в прошлом, я написал его на языке Ruby. Учиться на моих ошибках.

## Пожарные машины

Я живу на главной улице в центре Сиэтла, на полпути между пожарной частью и большинством городских пожаров (или мне так кажется). Соответственно, за эти годы у меня развилось праздное любопытство к пожарной части моего города.

К счастью (с точки зрения данных), существует сайт Realtime 911 (<http://www2.seattle.gov/fire/realtime911/getDatePubTab.asp>), на котором выложены все вызовы по пожарной тревоге вместе с номерами участвовавших пожарных машин.

И вот, чтобы побаловать свое любопытство, я собрал с их сайта многолетние данные о вызовах и выполнил анализ социальных сетей (<https://github.com/joelgrus/fire>) на примере пожарных машин. Среди прочего, потребовалось специально для пожарных машин изобрести понятие центральности, которое я назвал рейтингом пожарной машины.

## Футболки

У меня есть малолетняя дочь, и на протяжении всех ее дошкольных лет постоянным источником расстройств для меня было то, что большинство "девчачих" футболок были довольно скучными в то время, как большинство "мальчишечьих" выглядели забавными.

В частности, было ясно, что между футболками для девочек и футболками для мальчиков дошкольного возраста существует четкое различие. И поэтому я задался вопросом, можно ли обучить модель, которая распознавала бы эти различия?

Подсказка: ответ был положительным (<https://github.com/joelgrus/shirts>).

Для этого потребовалось скачать фотографии сотен футболок, сжать их все к единому размеру, преобразовать в векторы цвета пикселей и на основе логистической регрессии построить классификатор.

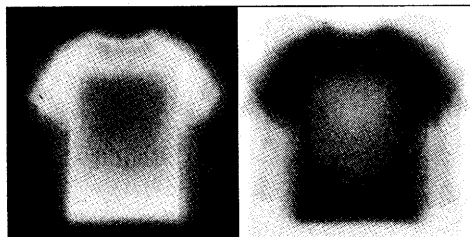


Рис. 25.2. Собственные футболки в соответствии с первой главной компонентой

На первом этапе исследовались цвета, присутствующие в каждой футболке. На втором были обнаружены первые 10 главных компонент векторов фотографий футболок, и каждая футболка была классифицирована при помощи их проекций на 10-мерное пространство, в котором расположились "собственные футболки" (рис. 25.2).

## **А вы?**

А что интересует вас? Какие вопросы не дают вам спать по ночам? Попробуйте отыскать набор данных (или извлечь их на веб-сайтах) и заняться наукой о данных.

Сообщите мне о своих находках! Шлите сообщения по электронному адресу [joelgrus@gmail.com](mailto:joelgrus@gmail.com) или найдите в Twitter по хештегу [@joelgrus](https://twitter.com/joelgrus).

---

# Предметный указатель

## A

Anaconda 17, 39

## B

Bokeh, проект 327

## G

GitHub 143

## I

IPython 324

## M

MapReduce 315

## N

NotQuiteABase, приложение 301

## P

pip 39

## S

Spyder 18

---

## A

Алгебра линейная 73

Алгоритм:

◇ "жадный" 231

◇ ID3 231

◇ PageRank 289

◇ Портера 198

Анализ:

◇ данных, интеллектуальный 170

◇ описательный 81

◇ регрессионный 215

◇ слов, морфологический 198

Анализатор, морфологический 198

## Б

База данных:

◇ NoSQL 313

◇ реляционная 301

Библиотека:

◇ BeautifulSoup 135, 137

◇ Bokeh 72

◇ D3.js 72, 327

◇ ggplot 72

◇ libsvm 225

◇ matplotlib 42, 63, 326

◇ NumPy 77, 80, 326

◇ pandas 92, 148, 168, 326

◇ requests 135

◇ scikit-learn 128, 168, 190, 198, 215, 225, 237, 263, 326

◇ SciPy 92, 263

◇ scipy.stats 105

◇ Scrapy 148

◇ seaborn 72, 327

◇ StatsModels 92

◇ Twython 145

Бизнес-модель 169

Бутстрап-агрегирование 236  
 Бутстрапирование 210  
 Бэггинг 236

## В

Вариация 85  
 Вектор 73  
 ◇ величина 76  
 ◇ длина 76  
 ◇ собственный 286  
 ◇ сумма квадратов 76  
 ◇ умножение на скаляр 75  
 Векторы:  
 ◇ вычитание 75  
 ◇ расстояние между векторами 77  
 ◇ сложение 74  
 Величина:  
 ◇ скалярная 73  
 ◇ случайная 97  
 ▫ биномиальная 103  
 ▫ равномерная 55  
 ▫ стандартная нормально  
 распределенная 100  
 Вероятность условная 94  
 Визуализация:  
 ◇ данных 63  
 ◇ интерактивная 63  
 Выражение регулярное 56

## Г

Гармоническое среднее 175  
 Генератор 54  
 ◇ (псевдо)случайных чисел 55  
 ◇ последовательностей 53, 78  
 ▫ с оператором for 54  
 Гипотеза:  
 ◇ альтернативная 106  
 ◇ нулевая 106, 114  
 ◇ статистическая 106  
 ▫ проверка 106  
 Гистограмма 65  
 Градиент 119  
 ◇ вычисление 120  
 Грамматика 269  
 Граф социальный 29  
 График линейный 68

## Д

Данные:  
 ◇ грязные 155  
 ◇ извлечение из веб-ресурсов 134  
 ◇ изменчивость 85  
 ◇ исследование 149  
 ◇ нормализация 161  
 ◇ поиск 327  
 ◇ сбор 129  
 ◇ управление 157  
 Дерево:  
 ◇ классификационное 227  
 ◇ принятия решений 226  
 ▫ создание 231  
 ◇ регрессионное 227  
 ◇ решающее 226  
 Диаграмма:  
 ◇ рассеивания *См. Диаграмма точечная*  
 ◇ столбчатая 65  
 ◇ точечная 70  
 Дисперсия 86  
 Дуга 79, 279  
 ◇ направленная 279  
 ◇ ненаправленная 279

## З

Значение:  
 ◇ максимальное 83  
 ◇ минимальное 83  
 ◇ среднее 75, 83  
 ▫ арифметическое 83

## И

Игра "20 вопросов" 226  
 Индекс 312  
 Интервал доверительный 111  
 Интерфейс программный 141  
 ◇ GitHub 143  
 Исключение 44  
 Истинность 51

## К

Капча 244  
 Каррирование *См. Применение аргументов  
 частичное*  
 Каузация 91

Квантиль 84  
 Класс 56  
 Классификатор наивный байесовский 178, 192  
 Кластер 250  
 ◇ объединение 258  
 ◇ расстояние между 258, 262  
 Кластеры, расстояние между 251  
 Ковариация 87  
 Конвейер обработки данных 130  
 Конструкция управляющая 50  
 Корреляция 88  
 ◇ в простой линейной регрессии 199  
 ◇ выбросы 88  
 ◇ ловушки 91  
 КORTEЖ 28, 46  
 Коэффициент:  
 ◇ детерминации 201  
 ◇ подобия, косинусный 294  
 Кураторство неавтоматическое 293

## Л

Лассо-регрессия 215  
 Лямбда-выражение 43

## М

Максимум 125  
 Математическое ожидание 98  
 Матрица 77  
 ◇ в виде векторов 78  
 ◇ единичная 78  
 ◇ корреляционная 153  
 ◇ разреженная 321  
 ◇ смежности 79  
 ◇ точечная 153  
 ◇ умножение на матрицу 284, 321  
 Медиана 83  
 Метод:  
 ◇ iteritems() 55  
 ◇ k средних 251  
 ◇ главных компонент 162  
 ◇ градиентного спуска 119  
 ▫ в множественной линейной регрессии 207  
 ▫ в простой линейной регрессии 202  
 ▫ пакетная минимизация (пример) 125  
 ▫ стохастического 126, 165  
 ◇ иерархической кластеризации, восходящий 257

◇ классификации на основе ближайших соседей 180  
 ◇ латентного размещения Дирихле 273  
 ◇ наименьших квадратов (МНК) 200  
 ◇ обратного распространения ошибки 243  
 ◇ опорных векторов 223  
 ◇ размножения выборок 210  
 ◇ случайных лесов 236  
 ◇ стохастического градиентного спуска, поиск оптимального значения коэффициента beta 208  
 Метрика:  
 ◇ F1 175  
 ◇ центральности 287  
 Минимум 120  
 Множество 50  
 Мода 85  
 Модели в машинном обучении 169  
 Моделирование:  
 ◇ прогнозное 170  
 ◇ тематическое 273  
 Модель:  
 ◇ без учителя 170  
 ◇ биграммная 267  
 ◇ компромисс между смещением и дисперсией 176  
 ◇ наименьших квадратов, допущения 206  
 ◇ параметризованная 170  
 ◇ правильность 173  
 ◇ распределенных вычислений MapReduce 315  
 ▫ анализ обновлений ленты новостей 319  
 ▫ базовый алгоритм 315  
 ▫ подсчет частотности слов 316  
 ▫ преимущества применения 317  
 ▫ сумматор 322  
 ▫ умножение матриц 321  
 ◇ с учителем 170  
 ◇ точность 174  
 ◇ языка,  $n$ -граммная 266  
 Модуль 41  
 ◇ random 55  
 Мощност проверки 109.

## Н

Наука о данных 324  
 ◇ решение задач (проекты автора) 328  
 Недообучение 171  
 Независимость 93

Нейрон 238  
 Новости хакера, сайт 328

## О

Облако слов 264  
 Обработка естественного языка 264  
 Обучение:  
 ◇ ансамблевое 237  
 ◇ без учителя 170, 250  
 ◇ контролируемое 250  
 ◇ машинное 170  
 □ переобучение и недообучение 171  
 □ правильность модели 173  
 ◇ неконтролируемое 250  
 ◇ с учителем 170, 250  
 Объект специальный None 51  
 Оператор:  
 ◇ break 51  
 ◇ continue 51  
 ◇ def 43  
 ◇ except 44  
 ◇ if 50  
 ◇ if-then-else трехместный 50  
 ◇ in 45, 50, 51  
 ◇ try 44  
 ◇ yield 54  
 Оператор SQL  
 ◇ CREATE TABLE 301  
 ◇ DELETE 304  
 ◇ FULL OUTER JOIN 311  
 ◇ GROUP BY 306  
 ◇ HAVING 306  
 ◇ INNER JOIN 310  
 ◇ INSERT 302  
 ◇ JOIN 309  
 ◇ LEFT JOIN 310  
 ◇ ORDER BY 308  
 ◇ RIGHT JOIN 311  
 ◇ SELECT 304  
 ◇ UPDATE 303  
 ◇ WHERE 304, 306  
 Операция арифметическая 42, 74  
 Оптимизация опыта взаимодействия 113  
 Отклонение стандартное 86  
 Отступ 40  
 Очередь 281  
 Ошибка:  
 ◇ в модели множественной линейной регрессии 207

◇ коэффициентов, стандартная 211  
 ◇ ложноположительная 108  
 ◇ при кластеризации 254  
 ◇ регрессии, случайная 199

## П

Пара "ключ-значение" 47, 55, 318  
 ◇ в словарях Python 315–317  
 Парадокс Симпсона 90  
 Переменная:  
 ◇ args 61  
 ◇ kwargs 61  
 ◇ спутывающая 90  
 ◇ фиктивная 178  
 Переобучение 171  
 Перцептрон 238  
 Плотность распределения вероятностей 98  
 Подгонка  $p$ -значения 113  
 Подзапрос 312  
 Поиск популярных тем 36  
 Полнота 175  
 Поправка на непрерывность 110  
 Правдоподобие 107, 117, 119, 212, 218, 219, 276  
 ◇ максимальное 203  
 Представление двоичное 79  
 Прецизионность 175  
 Признак 177  
 ◇ бинарный 178  
 ◇ извлечение 178  
 ◇ отбор 178  
 Присваивание множественное 47  
 Проблема "проклятия размерности" 186  
 Проверка односторонняя 109  
 Проект:  
 ◇ новости хакера 328  
 ◇ пожарные машины 329  
 ◇ собственные футболки 330  
 Произведение скалярное 75  
 Производная частная 119

## Р

Размах 85  
 ◇ интерквартильный 87  
 Размерность 151, 153, 161, 186  
 ◇ снижение 162  
 Распаковка аргументов 60  
 Распределение:  
 ◇ апостериорное 115

- ◇ априорное 115
- ◇ бета 115
- ◇ биномиальное 103
- ◇ вероятностей, непрерывное 98
- ◇ дискретное 98
- ◇ нормальное 100
- ◇ равномерное 98

Расстояние евклидово 77

Регрессия:

- ◇ гребневая 213
- ◇ линейная:
  - множественная 205, 207–209
  - простая 199
  - простая модель 199
- ◇ логистическая 218
  - задача предсказания оплаты аккаунтов 216
  - качество подбора модели 221
  - логистическая функция 218

Регуляризация 213

Рекомендации 292

- ◇ на основе популярности тем 293

## С

Сглаживание 193

Сеть 279

- ◇ нейронная:
  - искусственная 238
  - прямого распространения 240
- ◇ социальная, анализ 329

Символ пробельный 40

Скобка квадратная 40, 45–47

Словарь 47

◇ Counter 49

◇ defaultdict 48

◇ dict 282

Смещение 176

- ◇ дополнительные данные 177

События:

- ◇ зависимые 93
- ◇ независимые 93

Сортировка 52

Спам-фильтр 191

Список 45, 73

- ◇ ассоциативный *См. Словарь*

◇ метод sort 52

◇ объединение и разъединение 60

◇ распаковка 46

Статистика математическая 81

Строка:

- ◇ командная 129
- ◇ неформатированная 44
- ◇ символьная 44
- Сумматор 322
- Схема данных 301
- Сэмплирование по Гиббсу 271

## Т

Таблица:

- ◇ нормализованная 309
- ◇ сопряженности 174

Тенденция 68

Теорема:

- ◇ Байеса 96
- ◇ центральная предельная 103

Теория вероятностей 93

Тестирование A/B 113

Тип логический Boolean 51

Точка отсечения 36, 107, 110

Точность 175

Триграмма 268

Трюк ядерный 224

## У

Удаленность 181

Узел 279, 282, 288

◇ листовой 231

◇ решающий 231

◇ сети 79

Уровень значимости 108

## Ф

Файл:

- ◇ запись 131
- ◇ с разделителями 132
- ◇ текстовый 131
- ◇ чтение 131

Фильтрация коллаборативная

◇ на основе пользователя 294

◇ по схожести предметов 297

Функции math.erf 101

Функция 43

◇ %paste 41

◇ all 52

◇ any 52

◇ enumerate 59

◇ filter 59

- ◇ map 58
- ◇ partial 58
- ◇ range 53
- ◇ reduce 59
- ◇ xrange 54
- ◇ zip 60, 74
- ◇ анонимная 43
- ◇ компонентная 57
- ◇ предикативная 138
- ◇ распределения:
  - дифференциальная (ДФР) 98
  - интегральная (ИФР) 99
  - интегральная, обратная 101
  - кумулятивная 99
- ◇ сигмоидальная 240

## Ц

Центральность:

- ◇ по близости 283, 284
- ◇ по посредничеству 280
- ◇ по степени узлов 29, 280
- ◇ собственного вектора 284, 287

Цикл:

- ◇ for 41, 51, 53
- ◇ while 51

## Ч

Число:

- ◇ псевдослучайное 55
- ◇ случайное, генерация 55
- ◇ собственное 286

## Ш

Шкала данных 160

Шум 168, 171

## Э

Экземпляр класса 56

Энтропия 228

- ◇ разбиения 230

## Я

Язык

- ◇ R 327
- ◇ структурированных запросов (SQL) 301



# Data Science

## Наука о данных с нуля

Книга позволяет освоить науку о данных, начав «с чистого листа».

Она написана так, что способствует погружению в Data Science аналитика, фактически не обладающего глубокими знаниями в этой прикладной дисциплине.

При этом вы убедитесь, что описанные в книге программные библиотеки, платформы, модули и пакеты инструментов, предназначенные для работы в области науки о данных, великолепно справляются с задачами анализа данных.

А если у вас есть способности к математике и навыки программирования, то Джоэл Грас поможет вам почувствовать себя комфортно с математическим и статистическим аппаратом, лежащим в основе науки о данных, а также с приемами алгоритмизации, которые потребуются для работы в этой области.

В сегодняшнем хаотическом потоке данных скрыты ответы на многие волнующие человека вопросы. Книга познакомит с методологией, которая позволит правильно сформулировать эти вопросы и найти на них ответы.

### Вместе с Джоэлом Грасом и его книгой

- Пройдите интенсивный курс языка Python
- Изучите элементы линейной алгебры, математической статистики, теории вероятностей и их применение в науке о данных
- Займитесь сбором, очисткой, нормализацией и управлением данными
- Окунитесь в основы машинного обучения
- Познакомьтесь с различными математическими моделями и их реализацией по методу  $k$  ближайших соседей, наивной байесовской классификации, линейной и логистической регрессии, а также моделями на основе деревьев принятия решений, нейронных сетей и кластеризации
- Освойте работу с рекомендательными системами, приемы обработки естественного языка, методы анализа социальных сетей, технологии MapReduce и баз данных

**Джоэл Грас** работает инженером-программистом в компании Google. До этого занимался аналитической работой в нескольких стартапах. Активно участвует в неформальных мероприятиях специалистов в области науки о данных. Всегда доступен в Twitter по хештегу [@joelgrus](https://twitter.com/joelgrus).

*«Джоэл проведет для вас экскурсию по науке о данных. В результате вы перейдете от простого любопытства к глубокому пониманию насущных алгоритмов, которые должен знать любой аналитик данных».*

Роит Шивапрасад  
Специалист компании Amazon  
в области Data Science с 2014 г.

ISBN 978-5-9775-3758-2



**БХВ-ПЕТЕРБУРГ**

191036, Санкт-Петербург,  
Гончарная ул., 20  
Тел.: (812) 717-10-50,  
339-54-17, 339-54-28  
E-mail: mail@bhv.ru  
Internet: www.bhv.ru