

Кирилл Еременко

Работа с данными в любой сфере

Как выйти
на новый уровень,
используя аналитику



Москва
2019

*Моим родителям,
Александру и Елене Еременко,
которые научили меня самому важному
в жизни — быть хорошим человеком*

CONFIDENT DATA SKILLS

**MASTER THE FUNDAMENTALS OF WORKING WITH
DATA AND SUPERCHARGE YOUR CAREER**

KIRILL EREMENKO

УДК 004.6
ББК 65.291.213
Е70

Переводчик Д. Шалаева
Научный редактор З. Мамедьяров
Редактор Л. Любавина

Еременко К.

Е70 Работа с данными в любой сфере: Как выйти на новый уровень, используя аналитику / Кирилл Еременко ; Пер. с англ. — М. : Альпина Паблишер, 2019. — 303 с.

ISBN 978-5-9614-2582-6

Что общего у аналитика данных и Шерлока Холмса? Как у Netflix получилось создать 100%-ный хит — сериал «Карточный домик»? Ответ кроется в правильном использовании данных. Эта книга — практическое руководство и увлекательное путешествие в науку о данных, независимо от того, хотите ли вы использовать анализ данных в своей профессии, собираетесь ли стать аналитиком данных, или уже работаете в этой области. Ее автор, основатель образовательного онлайн-портала и консультант, Кирилл Еременко просто и понятно рассказывает об основных методах, алгоритмах и приемах, которые вам помогут на любом этапе: от сбора данных и их анализа до визуализации полученных результатов. Благодаря «Работе с данными в любой сфере» вы не только узнаете, как данные влияют на нашу жизнь (и как защитить свои данные), но и сможете расширить свои карьерные возможности.

УДК 004.6
ББК 65.291.213

Все права защищены. Никакая часть этой книги не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами, включая размещение в сети интернет и в корпоративных сетях, а также запись в память ЭВМ для частного или публичного использования, без письменного разрешения владельца авторских прав. По вопросу организации доступа к электронной библиотеке издательства обращайтесь по адресу mylib@alpina.ru

© Kirill Eremenko, 2018

This translation of *Confident Data Skills* is published by arrangement with Kogan Page.

© Издание на русском языке, перевод, оформление.

ООО «Альпина Паблишер», 2019

ISBN 978-5-9614-2582-6 (рус.)
ISBN 978-0-7494-8154-4 (англ.)

Содержание

Бонус для читателей	7
Введение.....	8

ЧАСТЬ ПЕРВАЯ «Что это?» Ключевые принципы

01	Определение данных.....	19
02	Как данные удовлетворяют наши потребности	39
03	Мышление, необходимое для эффективного анализа данных.....	55

ЧАСТЬ ВТОРАЯ «Когда и где я могу получить их?» Сбор и анализ данных

04	Сформулируйте вопрос.....	87
05	Подготовка данных.....	109
06	Анализ данных (часть I).....	133
07	Анализ данных (часть II).....	191

ЧАСТЬ ТРЕТЬЯ
«Как я могу это показать?»
Представление данных

08	Визуализация данных	223
09	Презентация данных	261
10	Ваша карьера в науке о данных	277
Благодарности		297
Литература		299

Бонус для читателей

Спасибо, что выбрали эту книгу. Вы сделали огромный шаг на пути в науку о данных.

Получите бесплатный доступ к моему курсу A-Z Data Science. Просто зайдите на сайт www.superdatascience.com/bookbonus и используйте пароль datarockstar.

Удачи в анализе данных!

Введение

«Наверное, вы всегда хотели стать аналитиком данных — с самого детства?»

Мне приятно, что меня об этом спрашивают. Да, я люблю свою работу. Я с большим удовольствием обучаю студентов основам науки о данных. И здорово, что люди, похоже, думают, что энтузиазм по отношению к данному предмету возник во мне еще в молодом возрасте. Но это абсолютно не соответствует действительности. Скажем честно, ни один ребенок не мечтает о том, чтобы стать ученым — аналитиком данных. Дети хотят быть космонавтами. Танцорами. Врачами. Пожарными. И если вы грезите о спасении жизней или о полетах в космическом пространстве, вы вряд ли остановите свой выбор на столь приземленном занятии.

Когда люди спрашивают меня, всегда ли я хотел построить карьеру в области науки о данных, я возвращаюсь к своему детству и вижу маленького русского мальчика, выросшего в Зимбабве. Запах тлеющих углей, брачные вопли африканских красных жаб, незабываемый уют зимнего вечера, кончики пальцев, переворачивающие страницу за страницей сборника историй для детей, — это фрагменты воспоминаний о множестве прекрасных вечеров, когда я слушал русские сказки, которые читала мне мама.

Моя мать хотела, чтобы я, мои братья и сестры любили Зимбабве, но она также заботилась о том, чтобы мы знали свои культурные корни. Она подумала, как наилучшим образом передать нам эту информацию, и решила, что самый действенный способ — сказки. Когда я в конце концов вернулся в Москву — в город, который едва помнил, — то почувствовал, что возвращаюсь домой, благодаря крупным информационным о России, вплетенным в затейливые сюжеты.

Такова сила повествования. И все множество услышанных сказок я хотел разбить на составляющие их компоненты. Мне нужно было

увидеть большую картину, но я хотел видеть ее сквозь призму маленьких деталей. Я был очарован каждой частью механизма, создающего что-то настолько прекрасное. Я интуитивно знал: для того чтобы самому рассказать хорошую историю, сначала нужно собрать эти маленькие единицы информации. Именно так сформировалось мое отношение к данным.

В сегодняшнюю цифровую эпоху данные используются для создания историй о том, кто мы такие, как мы себя представляем, что нам нравится и когда мы хотим чего-то. Для того, чтобы проложить тропинку с уникальными виртуальными следами. Машины теперь знают о нас больше, чем мы сами, благодаря всем доступным им данным. Они читают наши личные данные как сборник рассказов о нас. И в науке о данных замечательно то, что любая дисциплина сегодня записывает свои данные, а это значит, что, освоив профессию аналитика данных, мы также можем стать космонавтами, танцорами и врачами, о чем так сильно мечтали.

Мало кто знает, что работать с данными в конечном итоге означает быть рассказчиком, передающим информацию. Так же, как и структурные компоненты историй, проекты по анализу и обработке данных тоже организованы логически. В книге «Работа с данными в любой сфере» четко выделяются пять этапов, которые составляют то, что я называю процессом обработки и анализа данных. Это не единственный подход, который можно использовать, но он обеспечит нашему проекту связь с практикой и продвижение к логическому завершению. И он четко и ясно структурирован, что мне так нравилось в детстве.

И вот я решил рассказать историю данных...

Но я абсолютный новичок

Наука о данных фактически является одной из тех областей, которые извлекают выгоду из опыта других сфер. Я надеюсь, что многие мои читатели уже весьма преуспели в той или иной профессии. Хорошо. Вы *ничего* не потеряете, если обратитесь к науке о данных, работая в другой области. Отнюдь не вредно для начала разбираться в чем-то

еще. Это своего рода фундамент, который вам пригодится, чтобы стать хорошим аналитиком данных.

Начав работать в транснациональной консалтинговой компании Deloitte, я не знал ни одного из алгоритмов, которые мы рассмотрим в этой книге. Да никто от меня этого и не ожидал. Совсем немногие начали свою карьеру с науки о данных. Прочитав книгу, вы обнаружите, что те, кто добился успеха в этой сфере, даже не думали о ней, пока находились в начале своей карьеры. Итак, отбросьте страх перед цифровой неграмотностью — взяв эту книгу, вы сделали первый шаг на пути в мир науки о данных.

Эй, а где код?

Если вы, как и я, пролистываете книгу, прежде чем приступить к чтению, то, возможно, заметили, что вам не встретилось ни одной строки кода. Я слышу, как вы говорите: «Но это ведь книга о науке о данных, так что же происходит?» Наука о данных — чрезвычайно широкий предмет. «Работа с данными в любой сфере» погружает вас в тему и вдохновляет на размышления о том, как эта дисциплина может быть включена в вашу текущую или будущую деловую практику. Вы узнаете *методы* науки о данных — потому что ее «ингредиенты» (код) легко доступны онлайн. Если воспользоваться аналогией с приготовлением пищи, перед вами в меньшей степени просто книга рецептов и в большей — подробная информация об основных методах, используемых в науке о данных. Изучите их тщательно, и вы начнете интуитивно понимать, *почему* вам нужно применять определенные коды и методы, — гораздо более эффективный подход к обучению, чем просто предоставление строк кода для подключения к вашему проекту.

Как пользоваться этой книгой

Я написал эту книгу специально для того, чтобы вы могли обратиться к ней, где бы вы ни находились — в поезде, в ванне, в ожидании человека своей мечты. Читайте ее по частям или в один присест, по главам,

выбирая самое лучшее, выделяя нужное желтым маркером, наклейками. В начале каждой части вы найдете краткое введение, помогающее быстро определить, какая глава окажется для вас наиболее интересной. Часть первая более объемна, она дает общее представление о науке о данных. Вторая и третья части сосредоточены на процессах анализа и обработки данных, интуиции, стоящей за некоторыми из самых мощных на сегодняшний день аналитических моделей, и на том, как повысить ваши шансы на успех, совершая первые шаги в направлении цели.

Если вы новичок, то получите максимальную отдачу от книги, прочитав ее от корки до корки. Если вы знакомы с наукой о данных как с дисциплиной и хотите добраться до сути того, как применять ее методы, не стесняйтесь обратиться к главе, которая вам больше всего поможет.

ЧАСТЬ ПЕРВАЯ

«Что это?»

Ключевые принципы

Учитывая очевидно безграничный потенциал технических и прикладных наук и связанные с ними широкие возможности для умелых предпринимателей, некоторые могут спросить, почему они вообще должны заниматься наукой о данных — почему бы просто не изучить технологические принципы? В конце концов, технологии управляют миром и не выказывают никаких признаков сдачи позиций. Любой читатель, заботящийся о своей карьере, может подумать, что научиться разрабатывать новые технологии, несомненно, будет наилучшим способом двигаться вперед.

Легко расценивать технологии как фактор, который меняет мир, — они дали нам персональный компьютер, интернет, искусственные органы, беспилотные автомобили, глобальную систему позиционирования (GPS), — но мало кто думает о науке о данных как о движущей силе многих из этих изобретений. Вот почему вам стоит прочитать именно *эту* книгу, а не книгу о технологиях: вам нужно понять, как работает система, чтобы внести в нее изменения.

Мы не должны рассматривать данные только как скучных, но готовых помочь родителей, а технологии — как стильных подростков. Важность науки о данных не начинается и не заканчивается объяснением того, что технологии нуждаются в данных как одном из многих других функциональных элементов. Это было бы отрицанием прелести данных и множества интересных приложений, которые они

предлагают для работы и игры. Короче говоря, невозможно иметь одно без другого. Это означает, что, если у вас есть основа для науки о данных, перед вами будет открыта дверь к широкому кругу других областей, в которых нужен аналитик данных. Это делает науку о данных необычной и благоприятной областью исследований и практики.

В первой части приводится информация о вездесущности данных, а также о развитии и ключевых принципах науки о данных. Эти сведения полезны для начального погружения в предмет. Вы получите четкое представление о том, какое отношение данные имеют к вам, и задумаетесь не только о том, как данные могут непосредственно принести пользу вам и вашей компании, но и как вы можете в течение длительного времени использовать их в профессиональной и прочих сферах.

Начало пути

Глава 1 станет началом нашего путешествия в науку о данных. Сначала в ней будет продемонстрировано, насколько велики масштабы распространения данных и то, каким образом мы все вносим вклад в их производство в наш компьютерный век. Затем я расскажу, как люди собирают данные, работают с ними и, что очень важно, как данные можно использовать для поддержки большого количества проектов и методов внутри и вне самой дисциплины.

Мы установили, что проблемы с наукой о данных частично связаны не с ее относительной сложностью, а скорее с тем, что эта область знаний для многих по-прежнему покрыта туманом. Только когда мы точно понимаем, сколько данных имеется и как они собраны, мы можем начать рассматривать различные способы работы с ними. Мы достигли той точки в нашем технологическом развитии, когда информацию можно эффективно собирать и хранить на благо всех отраслей промышленности и научных дисциплин, о чем свидетельствует количество общедоступных баз данных и финансируемых правительством проектов по агрегированию данных культурными и политическими институтами. Вместе с тем сравнительно немногие знают, как получить доступ к данным и как их проанализировать. Если же люди

не осознают пользу данных для своей профессиональной деятельности, все красивые массивы данных только собирают пыль. В этой главе объясняется, почему наука о данных крайне важна *именно сейчас*, почему это не просто тенденция, которая скоро выйдет из моды, и почему вы должны рассмотреть возможность внедрения ее практик в качестве ключевого компонента решения ваших рабочих задач.

Наконец, в этой главе описывается, как стремительная траектория развития технологий не позволяет нам даже на время отвернуться от науки о данных. Каковы бы ни были представления о мире, к которому мы стремимся, невозможно остановить сбор данных, их обработку и использование. Тем не менее нельзя игнорировать тот факт, что сами по себе данные не касаются вопросов морали, и это обуславливает возможность их нечестного или неправильного использования. Те из вас, кто обеспокоен такого рода злоупотреблениями, могут принять участие в противостоянии им и вступить в дискуссию с глобальными институтами, которые занимаются проблемами, связанными с этикой данных — аспектом, который я нахожу настолько существенным, что отвел ему специальный подраздел в главе 3.

Будущее принадлежит данным

Все — каждый процесс, каждый датчик — скоро будет управляться данными. Это резко изменит способ ведения бизнеса. Я предсказываю, что через десять лет от каждого сотрудника любой организации в мире будет требоваться обладание определенным уровнем грамотности в сфере данных и умение работать с ними, получая на их основе некоторые идеи для повышения ценности бизнеса. Не такая уж дикая мысль, если учесть, что на момент публикации этой книги предполагается, что многие люди знают, как пользоваться цифровым кошельком Apple Pay, выведенным на рынок только в 2014 г.

Глава 2 — «Как данные удовлетворяют наши потребности» — наглядно демонстрирует, что данные являются эндемичными для каждого аспекта нашей жизни. Они управляют нами, накапливая силу в цифрах. Данные всегда играли важную роль в нашем существовании. Наша ДНК несет в себе основные данные о нас, и эти базовые формы

данных руководят нами: отвечают за то, как мы выглядим, за форму наших конечностей, за структуру нашего мозга и его способность обрабатывать информацию, а также за диапазон эмоций, которые мы испытываем. Мы — хранилища этих данных, шагающие флеш-накопители биохимической информации; вместе с данными нашего партнера мы передаем их нашим детям и «кодируем». Не интересоваться данными означает не интересоваться самыми фундаментальными принципами жизни.

В этой главе объясняется, как данные используются во многих областях, и для иллюстрации я использую примеры, которые непосредственно перекликаются с пирамидой потребностей Абрахама Маслоу, теорией, хорошо знакомой многим ученым и практикам в области бизнеса и управления. Если эта иерархия является для вас новинкой, не беспокойтесь — я объясню ее суть и то, как она применима к нам, в главе 2.

Приостановка развития

Последняя глава первой части покажет, как новички в науке о данных могут изменить свое мышление, чтобы погрузиться в нее, и поможет выявить те области, где уже сейчас возможно применить анализ данных. Многие достижения науки о данных основательно затронули другие сферы и поставили вопросы о будущем перед самыми разными специалистами и учеными. Если вы хотите развивать свою карьеру как аналитик данных, эта глава подскажет некоторые идеи для сфер, в которых вы, возможно, уже работаете.

В главе 3 я также представлю некоторые наиболее важные подходы, которые вы можете использовать, чтобы начать работу как практик. Наука о данных намного проще, чем многие другие научные дисциплины. Вам не нужно быть прирожденным ученым, чтобы овладеть принципами науки о данных. Что вам *действительно* необходимо — это умение придумывать различные способы извлекать пользу из данных тогда, когда дело касается бизнес-операций или личной мотивации. Ведь ученые — исследователи данных изучают *возможности* предоставленной информации. Вы можете удивиться, узнав, что у вас

уже есть некоторые навыки и опыт, которые вы можете использовать на своем пути к освоению этой дисциплины.

Разумеется, новичкам необходима разумная осторожность. Любой, кто использовал Excel, работал в офисной среде или изучал в университете предмет, имеющий научную составляющую, вероятно, уже встречался с данными. Но некоторые из методов использования данных, которые вы, возможно, усвоили, будут неэффективными, и приверженность тому, что вы уже знаете, может мешать вам изучить наиболее действенные способы использования массивов данных: мы обсудим это подробно во второй и третьей частях.

Несмотря на явный положительный эффект использования данных, важно не обольщаться. Поэтому в главе 3 рассматриваются и различные угрозы безопасности, которые данные могут представлять для своих пользователей, и то, как работают аналитики данных для решения текущих и потенциальных проблем. Этика данных является особенно привлекательной и заслуживающей внимания областью, поскольку она способна изменять и направлять будущие разработки в области науки о данных. Учитывая то, что мы знаем о сборе информации, этика данных — в той мере, в какой ее можно использовать в машинах и онлайн, — создает основу для общения людей и технологий. Когда вы прочитаете эту главу, подумайте о том, как каждая из областей может быть связана с тем, как вы работаете, и насколько полезны для вашего бизнеса дальнейшие инвестиции в эту сферу.

Подумайте о последнем фильме, который вы видели в кинотеатре. Как вы впервые узнали о нем? Возможно, вы кликнули на трейлер, когда YouTube рекомендовал его вам, или же ролик появился в качестве рекламы, прежде чем YouTube показал вам видео, которое вы действительно хотели посмотреть. Может быть, вы прочитали в социальной сети, что ваш друг хвалит картину, или в вашей новостной ленте появился увлекательный клип из фильма. Если вы любитель кино, сайт-агрегатор мог подобрать его для вас как фильм, который вам может понравиться. Вы, не исключено, нашли анонс фильма за пределами интернета — в своем любимом журнале либо же могли обратить внимание на афишу по дороге в кофейню, где лучше работает Wi-Fi.

Ни один из этих источников информации не был случайным. Звезды не просто сошлись для вас и фильма в нужный момент. Оставим идеалистические совпадения неожиданным экранным встречам. То, что привело вас в кино, было в меньшей степени желанием увидеть фильм и в гораздо большей — мощной смесью основанных на данных признаков, которые выделили вас в качестве вероятного зрителя, прежде чем вы сами поняли, что хотите посмотреть фильм.

Когда вы взаимодействовали с каждым из этих источников информации, вы оставили немного сведений о себе. Мы называем их выхлопными данными. Этот процесс не ограничивается вашим присутствием в онлайн и важен не только для создания социальных сетей. Независимо от того, используете ли вы социальные медиаплатформы, *нравится* вам это или нет, вы делитесь своими данными.

Так было всегда — мы просто научились лучше записывать и собирать их. Любое количество ваших ежедневных взаимодействий может способствовать этому «выхлопу». По дороге в лондонское метро вас запечатлевают камеры видеонаблюдения. Сев на поезд, вы добавляете

информацию в базу «Транспорт» статистических данных Лондона об использовании метро в час пик. Когда вы делаете закладки или выделяете страницы романа на своем устройстве для чтения Kindle, вы помогаете дистрибьюторам понять, что особенно понравилось читателю, и что они могли бы разместить в будущих маркетинговых материалах, и как глубоко читатели склонны погрузиться в роман, прежде чем остановиться.

Если вы наконец решите отказаться от испытаний в общественном транспорте и вместо этого поедете в супермаркет на автомобиле, выбранная вами скорость поможет GPS-сервисам показывать своим пользователям в режиме реального времени, насколько напряженный трафик в районе, и также позволит вашему автомобилю оценить, сколько еще времени остается, прежде чем вам стоит искать автозаправочную станцию.

И сегодня, когда вы выходите из этих точек соприкосновения, оставленные вами данные уже собраны и добавлены в «проект» о вас, который детализирует ваши интересы, действия и желания.

Но это только начало истории данных. Я расскажу вам о том, насколько действительно распространены данные. Вы узнаете основные понятия, которые пригодятся на пути к овладению наукой о данных, а также ключевые определения, инструменты и методы — они позволят вам применить навыки работы с данными к своей собственной деятельности. Эта книга расширит ваши горизонты, показывая, как наука о данных может использоваться в разных областях такими способами, которые прежде казались вам невозможными. Я опишу, как умение работать с данными может дать толчок вашей карьере и изменить ваш бизнес — будь то посредством идей, которыми вы впечатлите топ-менеджеров, или даже благодаря запуску стартапа.

Данные повсеместны

Прежде чем двигаться дальше, нужно уточнить, что подразумевается под данными. Когда люди размышляют о данных, они думают о том, как те активно собираются, хранятся в базах данных на непостижимых

корпоративных серверах и направляются на исследования. Но это устаревший взгляд. Сегодня данные гораздо более вездесущи*.

Все весьма просто: данные — это любая единица информации. Это побочный продукт любых действий, пронизывающих каждую часть нашей жизни не только в сфере интернета, но также в истории, географии и культуре. Наскальные изображения — данные. Музыкальный аккорд — данные. Скорость автомобиля, билет на футбольный матч, ответ на вопрос анкеты — все это данные. Книга — это тоже данные, как и глава в этой книге, как слово в главе, а также буква в слове. Им не нужно *быть собранными*, чтобы считаться данными. Их не нужно хранить в архиве организации, чтобы они считались данными. Значительная часть данных в мире, вероятно, пока не объединены в какой-либо базе данных.

Предположим, что в этом определении данных как единицы информации данные являются *осязаемым прошлым*. Весьма мудро, если задуматься. Данные — это прошлое, а прошлое — это данные. Запись всего, что можно отнести к данным, называется базой данных. И аналитики данных могут использовать их для лучшего понимания наших нынешних и будущих действий. Они применяют тот же принцип, что веками использовали историки: мы можем учиться на опыте истории. Мы можем учиться на наших успехах — и на наших ошибках, чтобы улучшить настоящее и будущее.

Единственный аспект данных, который в последние годы резко изменился, — наша способность собирать, организовывать, анализировать и визуализировать их в контекстах, которые ограничены только нашим воображением. Куда бы мы ни пошли, что бы мы ни покупали, какими бы ни были наши интересы, все эти данные собираются и систематизируются в тренды, которые помогают рекламодателям и маркетологам продвигать свои продукты к тем, кто в них заинтересован;

* Теперь вы, вероятно, привыкли к тому, что люди используют слово «данные» как множественную форму слова «данное» и что на самом деле правильно употреблять его с глаголами во множественном, а не в единственном числе. Вы можете упомянуть, что «данное» было впервые зафиксировано в 1645 г. как используемое в единственном числе Томасом Уркхартом и что только 60 лет спустя, в 1702-м, это слово стало использоваться как существительное во множественном числе. — *Здесь и далее, за исключением особо оговоренных случаев, прим. автора.*

которые показывают политические предпочтения членов правительства в соответствии с их происхождением или возрастом и которые помогают ученым создавать искусственный интеллект (ИИ), реагирующий не только на простые запросы, но и на сложные эмоции, этику и идеологию.

С учетом всех обстоятельств вы можете спросить: «Каковы же ограничения: что мы называем данными, а что — нет? Считаются ли фактические сведения о цикле цветения растения (количественные данные) такими же данными, как фиксация ученым культурного обычая, связанного с передачей умирающему родственнику букета цветов из родной страны (качественные данные)?» Ответ — да. Данные не дискриминируются. Не имеет значения, является ли рассматриваемая единица информации количественной или качественной. Качественные данные, возможно, были менее полезными в прошлом, когда не была достаточно сложной технология их обработки, но благодаря достижениям в алгоритмах, способных обрабатывать такие данные, этот недостаток быстро уходит в прошлое.

Говоря об ограничениях понятия «данные», еще раз вспомните, что данные — это прошлое. Вы не можете получать данные из будущего, если только вам не удалось создать машину времени. Но в то время как данные нельзя получить из будущего, с их помощью *можно* получить представление о грядущем и прогнозировать его. И именно способность данных восполнить пробелы в наших знаниях делает их настолько увлекательными.

Большие данные прекрасны

Теперь, когда мы разобрались, что такое данные, нужно по-другому взглянуть на то, где и как они фактически хранятся. Мы уже продемонстрировали наш широкомасштабный потенциал создания данных (это «выхлопные данные») и пояснили, что, трактуя их как единицу информации, мы создаем очень широкую концепцию того, что понимается под данными. Итак, если они где-то рядом, где все это *происходит*?

К настоящему времени вам, вероятно, доводилось слышать термин «большие данные». Проще говоря, большие данные — это название,

присвоенное массивам данных со столбцами и строками, которых настолько много, что они не могут быть обработаны обычным аппаратным и программным обеспечением в течение разумного промежутка времени. По этой причине сам термин является динамичным — то, что расценивалось как большие данные в 2015 г., уже не будет считаться большими данными в 2020-м, поскольку к тому времени будут разработаны технологии, легко справляющиеся с подобными объемами.

Три V

Чтобы можно было считать массив данных большими данными, должно быть выполнено хотя бы одно из трех условий:

- 1.** Объем данных — то есть размер массива данных (например, количество строк) — должен исчисляться миллиардами.
- 2.** Скорость, то есть то, как быстро собираются данные (например, потоковое видео в интернете), предполагает, что скорость генерируемых данных слишком высока для адекватной обработки с использованием обычных методов.
- 3.** Разнообразие. Это подразумевает либо разнородность типов информации, содержащейся в массиве данных, таком как текст, видео, аудио или файлы изображений (известные как неструктурированные данные), либо таблицы, содержащие значительное количество столбцов, которые представляют разные свойства данных.

Мы пользуемся большими данными в течение многих лет для всех видов дисциплин и гораздо дольше, чем вы могли бы ожидать, — просто до 1990-х гг. не было термина для их обозначения. Так что я вас шокирую: большие данные — это не большая новость. Это, конечно, не новая концепция. Многие, если не все, крупнейшие корпорации располагают огромными хранилищами данных об их клиентах, продуктах и услугах, которые собирались в течение длительного времени. Правительства хранят данные о людях, полученные в результате переписей и регистрации по месту проживания. Музеи хранят культурные

данные — от артефактов и сведений о коллекционере до выставочных архивов. Даже наши собственные тела хранят большие данные в виде генома (подробнее об этом в главе 3 «Мышление, необходимое для эффективного анализа данных»).

Короче говоря, если вы просто не в состоянии работать с данными, то можете назвать их большими данными. Когда ученые используют термин, они делают это не просто так. Он применяется, чтобы привлечь внимание к тому, что стандартных методов для анализа данных, о которых идет речь, недостаточно.

Почему такая суеда вокруг больших данных?

Вам может показаться странным, что мы только начали понимать, насколько значимыми могут быть данные. Но когда мы в прошлом собирали данные, единственное, что мешало нам превратить их во что-то полезное, было отсутствие технологий. В конце концов, важно не то, насколько огромны данные; важно, что вы с ними делаете. Любые данные, «большие» или иные, полезны, только если из них можно извлекать информацию, и до того, как была разработана соответствующая технология, чтобы помочь нам проанализировать и масштабировать эти данные, их полезность могла быть измерена только интеллектуальными возможностями человека, пытавшегося с ними совладать. Но для сортировки больших данных требуется более быстрый и мощный процессор, чем человеческий мозг. До технологических разработок XX в. данные хранились на бумаге, в архивах, библиотеках и хранилищах. Теперь почти все новые данные, которые мы собираем, хранятся в цифровом формате (и даже старые данные активно преобразуются в цифровые, о чем свидетельствует огромное количество ресурсов, сосредоточенных в таких цифровых собраниях, как Europeana Collections и Google Books).

Хранение и обработка данных

С изобретением компьютера появилась возможность автоматизации процесса хранения и обработки данных. Но большие массивы данных

увязли в первых машинах; ученым, работавшим с электронными массивами данных в 1950-х гг., приходилось ждать решения простой задачи несколько часов. Вскоре пришли к выводу, что для *правильной* обработки больших массивов данных — для установления связей между элементами и использования этих связей с целью получения точных и значимых прогнозов — нужно создавать информационные носители, которые могли бы управлять данными и справляться с их хранением. Разумеется, по мере совершенствования технологий, основанных на вычислениях, менялись и возможности компьютеров по хранению и обработке данных. И за последние 70 лет мы не только научились эффективно хранить информацию, но и смогли сделать эту информацию переносимой. Те же самые данные, которые в 1970-х гг. помещались только на 177 778 гибких дисках, к 2000-му могли поместиться на *одном флеш-накопителе*. Сегодня вы можете хранить все это и многое другое в облаке (хранилище с виртуализированной инфраструктурой, которая позволяет просматривать ваши личные файлы из любой точки мира)*. Когда вы в следующий раз обратитесь к личным документам, хранящимся в местной библиотеке, у вас на работе или просто в вашем мобильном устройстве, имейте в виду: вы фактически делаете то, что в 1970-х гг. потребовало бы использования более 100 000 дискет.

Когда новые технологии облегчили хранение данных, исследователи начали обращать внимание на то, как эти сохраненные данные могут быть использованы на практике. Как мы начали создавать порядок из хаоса? Вернемся к нашему предыдущему примеру — фильму, который вы в последний раз смотрели в кинотеатре. Вероятно, вы были выбраны, чтобы увидеть этот фильм, не проницательным маркетологом, сосредоточенно изучавшим соответствующие критерии, а умной машиной, которая изучила ваши «выхлопные данные» и сопоставила их с найденными ею демографическими сведениями о тех, кто увидел этот фильм и получил от него удовольствие. Это может казаться новинкой, но, как мы уже установили, данные и их (ручная)

* Облачные данные хранятся за пределами сайта и в основном перемещаются по подводным кабелям, которые укладываются на дно океана. Так что облако находится не в воздухе, как мы могли подумать, а под водой. Карту расположения этих кабелей можно найти на www.submarinecablemap.com.

обработка уже давно существуют. Некоторые из киностудий Голливуда еще в 1950-х гг. собирали данные о том, что конкретно — от актера до режиссера и жанра — хотела увидеть их аудитория, а потом преобразовывали эту информацию в демографические характеристики респондентов, включавшие в себя возраст, местожительство и пол. Даже в то время люди принимали способные изменить ход событий решения в соответствии с информацией, извлеченной из данных.

RKO Pictures

Почему RKO Pictures, одна из голливудских студий «Большой пятерки» в 1950-х гг., продолжала снимать Кэтрин Хепберн в своих фильмах? Потому что данные показывали, что это был беспроигрышный выбор, способный привлечь внимание людей и в конечном итоге заставить их пойти в кинотеатры.

Конечно, есть место и для интуиции. На первом кастинге режиссер Джордж Кьюкор нашел актрису странной, но также признал, что «она обладала огромным чувством, которое проявлялось даже в том, как она брала стакан. Я подумал, что она очень талантлива...» (Fowles, 1992). Вот пример интуиции.

Опираясь на данные о положительном восприятии Хепберн зрительской аудиторией, RKO позже смогла воспользоваться и интуитивными предположениями Кьюкора относительно таланта актрисы и превратить их в надежные прогнозы о том, что студия сможет и дальше зарабатывать свои миллионы.

Это произошло благодаря Джорджу Гэллапу — первому человеку, который рассказал руководителям Голливуда о возможности использовать данные для принятия решений и прогнозирования, включая подбор актеров на главные роли и определение того, в какой жанр наиболее целесообразно вкладывать деньги*.

* Гэллап был статистиком, впервые ставшим известным публике, когда разработал метод, с помощью которого он точно предсказал переизбрание Франклина Д. Рузвельта в 1936 г.

Чтобы помочь RKO сделать это, Гэллап собрал, объединил и проанализировал качественные и количественные данные, которые охватывали демографическую информацию о зрительской аудитории RKO и ее мнение о фильмах, выпускаемых киностудией. Собирая эти данные, Гэллап создал модель, которая в первый раз сегментировала аудиторию кинозрителей демографически, выделив тех, кто благоприятно реагировал на определенные жанры, — модель, которая может и будет использоваться в дальнейшем для выборки и анализа данных.

Разрекламированный как предсказатель, помогающий студиям разбогатеть, Гэллап быстро стал любимцем многих лидеров киноиндустрии США, проверяя по данным опросов и интервью отношение аудитории к персонажам различных лент, от мультиков Уолта Диснея до фильмов Орсона Уэллса*.

Своим успехом Гэллап был обязан только данным (возможно, его можно назвать первым высокооплачиваемым аналитиком данных в мире). Его усилия в области статистики привели к тому, что этот ресурс по-прежнему имеет ценность за пределами своего первоначального замысла, обладая потенциалом охвата *неструктурированных* данных: записанных интервью представителей аудитории, отражающих культурные и социальные ценности того времени. Возможно, Гэллап подозревал, что потенциал анализа данных может только расти.

Данные могут генерировать контент

Итак, что если после всех умных свидетельств, основанных на данных, вы возненавидели фильм, который недавно видели в кинотеатре? Ну, данные, возможно, не могут предсказать все, но они, безусловно, заставили вас занять место перед экраном. Иногда данные могут получить тройку за достижения, но они всегда получают отлично за усилия. И над первым уже работают. Вместо того чтобы привязывать

* Более подробно о новаторской работе Джорджа Гэллапа см. Ohmer (2012).

нужные демографические показатели аудитории к новому фильму или телевизионному сериалу, кинокомпания теперь находят способы использовать данные об аудитории, чтобы принимать обоснованные решения о предлагаемых публике развлечениях.

Но эта перемена влечет за собой необходимость в большем количестве данных. По этой причине сбор данных не прекращается, как только вы посмотрели выбранный для вас фильм; любые последующие комментарии, которые вы оставляете в социальных сетях или шлете по электронной почте, изменение ваших привычек просмотра фильмов в интернете генерируют о вас как о «кинозрителе» свежий массив данных, который учитывается в любых будущих рекомендациях, прежде чем наконец вы станете частью какой-либо демографической группы. Таким образом, по мере того как из подростка-эмо, интересующегося только демоническим пением, вы превращаетесь в любителя сложной сюрреалистической буффонады, которого все избегают на коктейльных вечеринках, ваши данные будут меняться вместе с вами и адаптироваться к этим колеблющимся предпочтениям.

В качестве примечания: еще более приятная новость состоит в том, что данные не отрицают ваших интересов. Если вы только *прикидываетесь* знатоком, но в действительности, как только опускаете шторы, до сих пор наслаждаетесь дрянными фильмами о зомби, ваши данные сохранят этот тайный вскормленный вами энтузиазм.

Конечно, обратная сторона медали в том, что ваши данные могут выдавать секреты, касающиеся ваших предпочтений. Имейте в виду, что данные — это запись действий, они не будут лгать на ваш счет. Некоторые даже тратят недюжинные усилия, чтобы скрыть свой «фактический» след на сайтах цифровых музыкальных сервисов, теша собственное тщеславие: они запускают альбом музыки, которая, по их мнению, служит в обществе признаком хорошего вкуса, но не слушают ее, так что их накопленные данные представят искаженную версию того, что им нравится. На мой взгляд, у этих людей слишком много свободного времени, но манипулирование данными тем не менее является важной темой, и со временем мы вернемся к ней.

Кейс: Netflix

Сериал «Карточный домик», выпущенный развлекательной компанией Netflix, впервые доказал индустрии, насколько сильны могут быть данные не только в том, что касается охвата нужной аудитории определенными разновидностями контента, но и в управлении фактическим *производством* контента.

Сериал — политическая драма — выпуска 2013 г. был первой проверкой того, как данные могут быть применены в производстве хитов. В преддверии создания «Карточного домика» Netflix собирала данные о своих пользователях. Полученные сведения о зрительских привычках позволили Netflix группировать свой видеоконтент в разнообразные и даже удивительные категории. Интерфейс скрывал от пользователей эти категории, но тем не менее они были использованы компанией, чтобы представить нужный фильм нужной аудитории.

Когда информация об этих подкатегориях появилась в интернете несколько лет назад, люди были ошеломлены. Чтобы вы могли получить представление о том, насколько точно действовала Netflix, вот некоторые варианты подкатегорий: «Захватывающие фильмы ужасов 1980-х», «Хорошее образование и воспитание с участием героев “Маппет-шоу”», «Драмы шоу-бизнеса», «Глуповатая независимая сатира», «Откровенные фильмы о реальной жизни», «Умные фильмы о заграничных войнах», «Бросающие в дрожь триллеры» и «Признанные критиками мрачные фильмы-экранизации». Таковы весьма специфические предпочтения зрителей. Но Netflix нашла значительную аудиторию для каждой из этих категорий и для многих других.

В конце концов исследователи данных в Netflix начали видеть совпадения в зрительских моделях их аудитории. Оказалось, что существует значительное число подписчиков Netflix, которые наслаждались и работой Кевина Спейси, и серьезными политическими драмами. Остальное — перезапуск оригинального «Карточного домика» 1990-х гг. с Кевином Спейси в главной роли — это история (или это данные?).

Оседлав волну успеха

Netflix оказалась права, высоко оценив возможности данных: сериал «Карточный домик» был отмечен наградами и получил высокие оценки критиков. Поэтому неудивительно, что многие конкуренты Netflix попытались скопировать эту выигрышную модель. Хейделин де Понтевес, предприниматель в области данных и мой бизнес-партнер, работал на конкурента Netflix в целях создания подобной системы.

«Мы знали, что у Netflix уже есть мощная система рекомендаций, и поэтому от нас как разработчиков баз данных и операционных систем требовалось не создать то же самое для нашей компании, а найти, где можно добиться разницы. Мы поняли, что для разработки действительно интересной системы нам нужно сделать больше чем просто инструмент для рекомендаций фильмов, соответствующих определенным демографическим сегментам. Мы также хотели создать алгоритм, позволяющий предлагать фильмы, которые могли бы вывести пользователей из их зоны комфорта, но в то же время доставить им удовольствие. Мы действительно стремились к тому, чтобы появился некий элемент неожиданности».

(Де Понтевес, 2017 г.)

Хейделин понимал, что для достижения этой цели потребуется сложная система, способная проникнуть в головы пользователей и понять их предпочтения лучше, чем те сами понимали это. Он достиг цели, извлекая все имевшиеся у компании данные по клиентам и применяя правильное сочетание моделей, чтобы найти связи между зрительскими привычками. Помните, что этот подход почти такой же, как был у Джорджа Гэллага многие годы назад; благодаря доступным технологиям и воображению аналитика данных мы теперь можем получить доступ к данным гораздо более хитроумным (и автоматизированным) способом.

Использование данных

Некоторые могут посетовать, что такой подход к использованию данных для творческого контента фактически убивает творчество. На это я бы ответил им, что данные всего лишь следуют за тем, чего хотят люди. Для любой отрасли желательно показать нужной аудитории в нужное время и в нужном месте соответствующий контент, чтобы побудить клиентов покупать их услуги. Таким образом, данные сделали индустрию более демократичной, потому что, хотя машины могут начать влиять на наши предпочтения в покупках, мы по-прежнему сохраняем самую ценную информацию: человеческое желание. Машины не говорят нам, чего мы хотим; они создают для нас связи, о которых мы, возможно, не знали.

Данные не приказывают людям идти и смотреть фильмы о супергероях и не смотреть французские сюрреалистические фильмы; они прислушиваются к тому, чего люди хотят и от чего получают удовольствие*. Если вы считаете, что существует проблема удушения творчества, то это не вина данных — это вина нашего общества. Я не устану подчеркивать, что данные *являются* прошлым. Это всего лишь запись информации. Если вы *хотите* видеть больше французских сюрреалистических фильмов, то просто идите и смотрите их — и убедитесь, что после просмотра вы о них говорите**. Может показаться, что вы просто добавляете шума в интернете, но этот «шум» быстро обрабатывается и становится доступным для использования повсюду. Благодаря данным в нынешнюю эпоху наши голоса действительно могут быть услышаны и иметь реальную власть — так почему бы не воспользоваться этим?

Кроме того, модели для использования данных еще несовершенны. В случае с медиаиндустрией другие корпорации приняли концепцию Netflix, и некоторые могут отметить, что одни преуспели больше, а другие — меньше. Но опять же, в этом нет заслуги данных, это творческий вклад людей. В конце концов, *именно здесь* находится нынешний предел нашей способности использовать данные для создания контента. Наверное, мы сможем оценить *вероятное* число людей, заинтересованных в концепции, но на карту поставлено гораздо больше, так как конечный успех любой формы развлечений будет обусловлен талантом ее создателя. Пусть это станет предупреждением для писателей и режиссеров, которые надеются получить легкие результаты, полагаясь исключительно на данные: базы

* Пример того, какие проблемы и возможности связаны с аналитикой данных в киноиндустрии, см. у Mishra and Sharma (2016), в докладе которых анализируется кинопроизводство и продюсирование в Индии.

** Естественно, на пути этого подхода есть препятствия. Вы не сможете победить миллионы поклонников супергероев в Китае, которые в значительной степени отвечают за то, что Голливуд продолжает наращивать выпуск фильмов о мужчинах (и женщинах) в колготках, спасающих мир от зла. Вопросы о том, как данные влияют на творчество, возможно, выходят за рамки этой книги, но я бы сказал, что всегда существовало и всегда будет существовать пространство для творчества, даже в мире, управляемом данными. Мы не становимся тупее; мы просто делаем промышленность более эффективной.

данных, которые показывают успех фильмов разных жанров, могут быть полезным руководством для последующих действий, но будут оставаться только руководством, поскольку результат работы зависит от таланта человека.

Почему данные важны сейчас

Многие уже в курсе того, что технологии в будущем могут существенно повлиять на рабочие места. Если вы чувствуете себя достаточно смелым, введите в поисковую строку Google «технологическое воздействие на рабочие места» /«technological impact on jobs» — и вы увидите, что несметное количество статей посвящено вероятности автоматизации в сфере вашей деятельности*. Хотя эта информация подкреплена данными, я бы сказал, что, возможно, мнение исследователей в некоторой степени субъективно, если принять во внимание задачи, которые необходимо выполнять на конкретных рабочих местах. Так, я бы, конечно, не рекомендовал учиться на спортивного арбитра по той причине, что эта работа зависит от *данных* об игре, — машины неизбежно будут поставлять более точные данные, чтобы подтвердить или опровергнуть любые заявления соперников. Судья может быть данью традиции, которая делает опыт более личностным или захватывающим *прямо сейчас*, но, на мой взгляд, ностальгия, связанная с профессией, не означает, что она будет востребована вечно.

Даже после того, как выяснилось, насколько всепоглощающими являются данные, некоторые все еще могут надеяться на то, что наука о данных не повлияет на их бизнес в ближайшее время. В конце концов, нужно время, чтобы что-то произошло. Но думать таким образом было бы большой ошибкой, потому что это отрицало бы принцип закона Мура.

* Опасения по поводу технологической безработицы не новы — Джон Мейнард Кейнс писал об этом в 1930-х гг.: «Мы страдаем от новой болезни, названия которой некоторые читатели, возможно, еще не слышали, но о которой они многое услышат в ближайшие годы, а именно — о технологической безработице» (Кейнс, 1963).

Закон Мура

Закон Мура — это закон прогнозирования. Предложенный соучредителем Intel Гордоном Муром в 1965 г., он в первую очередь касался ожидаемого со временем увеличения числа транзисторов (устройств, используемых для управления электрическим током) на квадратный дюйм в интегральных схемах (например, компьютерных микросхемах, микропроцессорах, материнских платах). Было замечено, что число этих транзисторов примерно удваивается каждые два года, и закон утверждал, что тенденция будет продолжаться. На сегодняшний день это подтвердилось*.

В восприятии непрофессионала это означает, что, если вы пойдете в свой местный компьютерный магазин сегодня и купите компьютер за £1000, а через два года приобретете еще один тоже за £1000 в том же магазине, вторая машина будет в два раза мощнее, хотя она стоит столько же.

Многие применили этот закон к растущему как грибы количеству достижений в области науки о данных. Она является одной из самых быстроразвивающихся академических дисциплин, и занимающиеся ею профессионалы используют все более изощренные способы, чтобы найти новые средства для сбора данных, построения экономичных систем их хранения и разработки алгоритмов, которые превращают все эти порции больших данных в ценные идеи. Доводилось ли вам когда-либо чувствовать, что технологии движутся вперед так быстро, что вы не успеваете идти в ногу со временем? Тогда подумайте об аналитиках данных. Они играют в салочки с технологией, которая *еще даже не изобретена*.

Кейс: Siri

В качестве примера рассмотрим развитие технологии распознавания речи. Создатели Siri Даг Киттлаус, Адам Чейер и Том Грубер разработали умного личного

* Относительно транзисторной инфраструктуры у закона Мура есть ограничения. При размере около 1 нм свойства полупроводникового материала нарушаются такими квантовыми эффектами, как квантовое туннелирование. Кроме того, дальнейшее развитие инфраструктуры потребует альтернативы кремнию, который сейчас используется в качестве основного материала. — *Прим. науч. ред.*

помощника задолго до того, как технология стала достаточно зрелой, чтобы можно было реализовать идеи и вывести их на рынок. Авторы Siri создали инструменты и алгоритмы для работы с имевшимися у них данными, чтобы поддерживать технологию распознавания речи, которая тогда еще не была изобретена.

Однако они знали, что, хотя было невозможно использовать программное обеспечение с имевшейся в то время технологией, в конечном итоге запуск Siri *станет возможным*, нужно лишь подождать, пока технология выкристаллизуется. Короче говоря, они уловили технологические тенденции.

Концепцией, которую создатели Siri использовали для своих прогнозов, служил закон Мура. И это невероятно важно для науки о данных. Закон Мура применяется к многим технологическим процессам и является необходимым правилом при рассмотрении и принятии деловых решений и реализации проектов; мы вернемся к его обсуждению в главе 3 «Мышление, необходимое для эффективного анализа данных».

Беспокойство ни к чему не приводит

Голливуд и индустрия развлечений в целом долгое время придерживались мрачной идеи, что использование данных и связанные с ними злоупотребления угрожают человечеству. Нам стоит задуматься над этой не предвещающей ничего хорошего фразой из фильма «2001: Космическая одиссея»: «Открой дверь модульного отсека, ЭАЛ», где ЭАЛ — технология искусственного интеллекта (ИИ) космического корабля — настолько усовершенствован, что решает не подчиняться команде человека и действовать согласно своим (превосходящим) суждениям. «Из машины», «Она», «Бегущий по лезвию», «Призрак в доспехах» — все эти фильмы посвящены воображаемым проблемам, с которыми могут столкнуться люди, когда технологии начнут развивать собственное сознание и предвидеть наши действия.

Но есть, с моей точки зрения, еще одна область, где злонамеренное применение данных — имеющее значительно больше общего

с злоупотреблениями *людей*, чем роботов, — гораздо более вероятно и неотвратно. Речь идет о конфиденциальности. С вопросами конфиденциальности связаны многие наши взаимодействия в интернете. Люди могут оставаться анонимными, но информация о них всегда будет где-то собираться — и использоваться. Даже если эти данные лишены характерных индикаторов, отсылающих к тому или иному индивидууму, некоторые могут спросить: «Правильно ли, что такие данные вообще собирают?»

Ваш онлайн-след

Читатели, которые пользовались интернетом в 1990-х гг., знакомы со словом «аватар» — довольно безобидное изображение, которое мы выбирали для представления себя на онлайн-форумах. Сегодня термин «аватар» используется для описания чего-то гораздо более широкого. Теперь он означает нашего неосязаемого двойника в виртуальном мире, массив данных о нас, составленный на основе наших заданных поисков, выбора и покупок, которые мы делаем в интернете, и всего, что мы публикуем в Сети, от текста до изображений. Такие данные являются потенциальным золотым дном, неиссякаемым источником информации для кредитных агентств и компаний-агрегаторов, которые затем могут использовать эти сведения для продажи другим.

Ввиду развития науки о данных встают вопросы этики и безопасности, касающиеся проникаемости, искажения и захвата данных (а этика — это область, которую мы рассмотрим в главе 5 «Подготовка данных»). У нас есть очень веские основания беспокоиться о доступах, которые открывает наука о данных, и о том, что она не делает различий в том, кто — или что — обращается к этой информации. Хотя переход от бумажного к цифровому документообороту позитивно сказался на практике ведения дел в компаниях, данные все еще могут пропадать или приходить в негодность, а также на них может существенно повлиять человек (это касается неверной информации, потери баз данных и шпионажа), что будет иметь разрушительные последствия.

Кейс: The Heartbleed Bug

На мой взгляд, Heartbleed Bug* представляет собой самое радикальное нарушение конфиденциальности в мире на сегодняшний день. Ошибка в программе позволила хакерам применить уязвимость в исходном коде, используемом в интернете, и украсть защищенные иным образом данные, отправленные через безопасные соединения Secure Sockets Layer (SSL). Эта лазейка предоставила доступ к конфиденциальной информации о торговых сайтах за много лет, прежде чем стало известно о ее масштабах.

В 2014 г. группа безопасности Google обнаружила эту проблему в исходном коде SSL во время регулярного критического просмотра своих сервисов. Оказалось, что около 800 000 веб-сайтов во всем мире имели эту ошибку в своем исходном коде, что обеспечивало доступ к их информации ворам и хакерам, знавшим об этой уязвимости. Но в течение двух лет ошибка оставалась незамеченной, что позволило украсть потенциально бесчисленное количество данных. По иронии, как сайты с поддержкой SSL (те, что начинаются с «https») они должны быть более безопасными, чем те, у которых обычные URL-адреса «http».

Даже если проигнорировать распространенное в то время мнение, что ошибка сохранялась с ведома правительственных или фиктивных организаций, факт остается фактом: Heartbleed Bug представлял собой фундаментальное нарушение конфиденциальности.

Не контролируйте — просвещайте!

Неудобная истина, касающаяся науки о данных и любой отрасли, где напрямую задействованы деньги, заключается в том, что по мере роста интереса к данной дисциплине возрастает интерес к наиболее гнусным средствам вмешательства в ее внутренние процессы. Некоторые могут счесть это достаточным основанием для прекращения сбора и использования данных. Но я вижу это по-другому и сделал бы ставку

* Ошибка в программном обеспечении OpenSSL, которая позволяет несанкционированно читать оперативную память. Вызывает двустороннюю уязвимость: не только вы можете читать данные с уязвимого сервера, но и злоумышленник оказывается способен получить доступ к вашей оперативной памяти, если у вас поврежденная версия OpenSSL. — *Прим. науч. ред.*

на то, что многие другие ученые — специалисты в области данных — чувствуют то же самое: вместо того чтобы контролировать и ограничивать, нужно воспитывать людей. Мы должны сообщить нашим детям, что их деятельность в интернете приведет к появлению аватара, который может быть использован в их пользу — или против них. Мы должны убедиться, что люди в целом лучше разбираются в том, как используют их данные и зачем.

Таков мир, в котором мы сейчас живем. Нам будет намного легче избавиться от этой эмоциональной привязанности, чем сопротивляться. В конце концов, сегодня на сцену выходит молодое поколение и рекламируются новые потребительские компании. Это подтверждается тем, что многие компании (от Amazon до Outfittery*) работают в интернете. Сейчас потребители готовы предоставить свою личную информацию в обмен на лучшую адаптацию продуктов и услуг к их потребностям. Посмотрите на Instagram или Twitter, и вы убедитесь, что передача личной информации в интернете — в самых разных областях — может восприниматься как вторая натура поколения миллениалов. Если вы не планируете жить вне Сети на лоне дикой природы и говорить только с птицами, кибербезопасность — просто еще один риск нынешней действительности. Борьба с этой угрозой будет так же бесполезна, как усилия луддитов в XIX в.: сколь яростно ни сопротивлялись они автоматизации производства, это мало что изменило в долгосрочной перспективе.

Намного менее вероятно то, что мы откажемся от услуг, которые уже интегрированы в нашу жизнь и считаются само собой разумеющимися, — прежде всего потому, что сейчас *мы нуждаемся* в них. Когда-то эти услуги были роскошью, но технологии быстро превратились в основную потребность, определяющую то, как мы живем и работаем. Технологии развивались, и впредь нам нужно использовать данные.

Когда одни события в мире быстро сменяются другими и есть возможность наблюдать за ними в режиме реального времени с помощью социальных сетей, настойчиво выплескивающих информацию, или

* Базирующаяся в Берлине компания, торгующая мужской одеждой. Продает коробки с индивидуально подобранными для каждого клиента товарами.

новостных сайтов, которые постоянно, в любое время суток, обновляют свои страницы, публикуя новые сведения о происходящем, — в этих условиях можно испытывать чувство подавленности. Лавина данных наступает со всех сторон, и нет способа ее остановить. Вы не можете заткнуть вулкан пробкой и ожидать, что он не взорвется.

Однако мы *можем* управлять данными и анализировать их. Вероятно, вы слышали о «кураторах контента» и «сайтах-агрегаторах», таких как Feedly, через которые можно отбирать и классифицировать новости из интересующих вас блогов и сайтов. Люди и компании работают над организацией важных для них самих или их подписчиков данных. Эти попытки управления информацией должны обеспечить нам комфорт, и они представляют собой одну из множества опций обработки данных. По мере совершенствования технологий, которые помогут нам управлять данными и анализировать их, мы примем это как неотъемлемую часть нашего существования в компьютерную эпоху. Поэтому отбросьте свои сомнения и давайте вместо этого сосредоточимся на возможностях данных и на том, как они могут улучшить нашу жизнь.

Как данные удовлетворяют наши потребности

В науке о данных не так уж много таинственного — она, в конце концов, полностью вписана в современные реалии. И все же преобладает неверное представление, будто данные сложны и даже непостижимы. К сожалению, многие сегодня либо охотно отказываются видеть, как широко применяется наука о данных, либо намеренно отвергают ее как нечто недоступное или неприменимое к их работе. Наука о данных как дисциплина *предполагает* что-то весьма замысловатое. Это похоже на то, чем люди занимаются в маленьких кабинетах без окон, сгорбившись над своими столами.

Такой взгляд совершенно неверен.

В этой главе мы точно узнаем, насколько данные вездесущи, как широко они генерируются и собираются и почему наука о данных никогда не может считаться причудой.

Проникновение данных

Чтобы проиллюстрировать, насколько важны данные для всех аспектов нашей жизни — что это необходимость, а не роскошь, я буду использовать пирамиду потребностей Маслоу, которая, я уверен, знакома многим бизнес-практикам. В литературе по бизнес-психологии о ней написано очень много. Я считаю, что эта модель на удивление хорошо сочетается с распространенностью и преимуществами данных*.

* Это не даст нам исчерпывающих сведений о том, как и где наука о данных используется в нашей жизни, поскольку пирамида Маслоу принижает неосновные человеческие

Иерархия потребностей была разработана Абрахамом Маслоу в 1943 г. для отображения сложной мотивации, обуславливающей поведение людей. Иерархия представлена в форме пирамиды, которая в последовательности снизу вверх включает в себя потребности — от наиболее к наименее фундаментальным (рис. 2.1). Короче говоря, иерархия организована таким образом, что потребности, находящиеся на самом нижнем уровне пирамиды, должны быть удовлетворены до того, как у индивидуума, о котором идет речь, появится мотивация для удовлетворения потребностей более высоких уровней*.



Рис. 2.1. Пирамида потребностей Маслоу

потребности. Такие области, как, например, военная оборона и освоение космического пространства, сюда не будут включены, поскольку они не являются основными потребностями человека.

* Я использую пирамиду потребностей Маслоу в качестве примера для описания всеобъемлющей силы данных, но, если вы хотите узнать больше о том, как эта иерархия может быть применена в бизнесе, см. Conley (2007).

Наука о данных и физиология

В основе иерархии Маслоу лежат физиологические факторы — основные потребности людей для простого выживания. Как данные могут поспособствовать лучшему удовлетворению этих основных потребностей?

Давайте возьмем в качестве примера воздух, которым мы дышим. Загрязнение воздуха — один из наиболее серьезных поводов для глобального беспокойства со времен промышленной революции конца XVIII и начала XIX в. Мы могли бы считать смог феноменом прошлого — так, в 1950-х гг. выбросы, образовавшиеся при сгорании угля, регулярно окутывали Лондон. Но смесь дыма, тумана и пыли остается большой проблемой во многих городах по всему миру, от Китая до Бразилии.

Любые технологии, предназначенные для уменьшения загрязнения воздуха в городах, зависят от данных: чтобы улучшить состояние воздуха, его состав необходимо сначала контролировать.

Кейс: экологические данные и «Зеленый горизонт»

Программа «Зеленый горизонт» (Green Horizon) была запущена компанией IBM в 2014 г. в связи с необходимостью отреагировать на ужасное качество воздуха в Китае путем «преобразования его национальных энергетических систем и поддержки потребностей в устойчивой урбанизации» (IBM, 2017a). «Зеленый горизонт»^{*} использует данные 12 глобальных исследовательских лабораторий и применяет когнитивные модели к собранным данным, чтобы предоставить информацию, связанную с главной целью проекта — сокращением загрязнения. Данные необходимы для мониторинга колебаний загрязнения воздуха в отдельных районах, а также для того, чтобы ученые могли проанализировать различные факторы, которые прямо или косвенно влияют на качество, температуру и состояние воздуха, и начать улучшать физическую среду в Китае.

Огромное преимущество этих проектов заключается в том, что экологические данные чаще всего являются общедоступными и в глобальном масштабе. Это

^{*} Программа использует интернет вещей и ИИ, чтобы предсказывать уровень загрязнения воздуха. — *Прим. науч. ред.*

означает, что технологические разработки, направленные на борьбу с загрязнением воздуха, могут быстро развиваться. Наличие доступа к важным массивам данных, связанных с удовлетворением наших самых основных потребностей, необходимо для понимания того, как имеющиеся технологии могут работать лучше. Вот почему у нас теперь есть специальные стеклянные панели, которые могут быть установлены в зданиях, чтобы окна могли «дышать», очищая воздух внутри помещения и тем самым защищая находящихся там людей. Вот почему у нас есть фильтры, которые могут быть использованы на фабриках в целях уменьшения вредных выбросов и защиты местных жителей от отравления.

Возобновляемые продовольственные ресурсы

Еда еще один пример того, как данные связаны с самыми основными потребностями человека (физиологические факторы в пирамиде Маслоу). Для некоторых это может показаться научной фантастикой, но уже в течение многих лет еда выращивается в лабораториях, а использование искусственного мяса становится все более актуальным феноменом. Memphis Meats, стартап в Кремниевой долине, который с момента своего создания разработал разные виды искусственного мяса, от говядины до домашней птицы, — всего лишь один из подобных институтов.

Поскольку это все еще некая «серая» область для регулирующих органов, религии и науки, искусственное мясо вызвало и похвалы, и гнев мирового сообщества (Devitt, 2017). Но нравится нам это или нет, искусственное мясо в недалеком будущем может стать заменой того, что мы едим. Резко сократив потребление воды и выбросы углерода, оно станет экологически безопасным решением в условиях, когда сельское хозяйство негативно влияет на мир природы. И данные, которые мы собираем для производства такого мяса, в конечном итоге выйдут за рамки исследования ДНК. Поскольку пищевые технологии становятся все более обыденными, дополнительные потребительские данные будут использоваться для других целей, таких как определение оптимальных способов приготовления искусственного мяса, — это позволит не только сделать мясо вкуснее, но и, что особенно важно для производящих компаний, повысить его продаваемость.

Наука о данных и безопасность

Как только физиологические потребности оказываются удовлетворены, приоритетом, согласно пирамиде Маслоу, становится безопасность (физическая, финансовая, личная). Таким образом, безопасность — это уровень, который в значительной степени включает в себя личное здоровье и благополучие, а медицина — одна из тех областей, для которых наука о данных особенно важна. В медицинской промышленности наука о данных радикально меняет инструменты для диагностики и лечения болезней. Все медицинские эксперименты проводятся с опорой на данные участников, и эти собранные данные могут использоваться для уточнения диагноза, подбора разных практических подходов и создания новых продуктов. Чтобы выявить сложные и редкие заболевания, практикующие медики должны владеть информацией о различных их проявлениях и симптомах — это поможет избежать ошибки при постановке диагноза, найти корень проблемы и эффективно ее решать. Когда недуг усугубляется и требует безотлагательного вмешательства врачей, течение болезни может не контролироваться на протяжении недель и месяцев, которые уходят на то, чтобы пациенты записались на прием к нужному специалисту.

От ученых — аналитиков данных требуется разработать передовые алгоритмы и обучить им машины для получения наиболее точных данных. На основе этих данных могут быть спрогнозированы необычные ситуации. Более того, собранные данные не зависят от благополучия научного сотрудника, работающего с ними (извините). Как только специалисты-медики выходят на пенсию, вместе с ними уходят их специфические знания. Когда аналитики данных уходят на заслуженный отдых, алгоритмы, которые они оставили, или собранные ими данные могут использоваться как основа для расширения существующих знаний. Наука о данных всегда опирается на то, что осталось, на информацию о нашем прошлом.

Именно эта способность позволяет столь эффективно использовать плоды науки о данных в медицине: пока данные сохраняются, накопленные знания не будут зависеть от отдельных людей.

Кейс: диагностика с помощью SkinVision

На рынке существует множество цифровых приложений, которые собирают данные по различным темам, от звезд в ночном небе до веснушек на вашей коже.

SkinVision — это приложение для мобильных устройств, помогающее тестировать родинки пользователей, чтобы выявить рак кожи. Используя агрегированные пользовательские данные, алгоритм SkinVision может определить вероятность появления у пользователя родинки с злокачественными симптомами. Это действительно очень просто: с помощью приложения вы делаете фото вашей кожи, SkinVision его регистрирует и проанализирует — а потом вы получите рекомендацию относительно следующих шагов, которые вы можете предпринять вместе с врачом.

Не стоит думать, что ставить диагноз с помощью мобильного устройства легкомысленно. По мере того как будет собрано все больше и больше сведений о болезни, базы данных о ее причинах и последствиях увеличатся и станут определять диагноз намного лучше, чем это делает опытный хирург. Чем больше людей используют цифровое приложение подобное SkinVision, чтобы узнать свой диагноз, тем выше вероятность, что технология сможет отличить доброкачественную родинку от злокачественной, потому что у нее будет большой массив данных, с помощью которых можно перекрестно изучить пользовательские данные — представленные изображения. Подумайте, что бы вы предпочли: получить диагноз от человека, которому довелось рассмотреть 1000 отдельных случаев, или от машины, которая накопила информацию о миллионе отдельных случаев?

Объем знаний

Отнюдь не только цифровые приложения прокладывают путь медицине, основанной на данных. Суперкомпьютер IBM Watson, по словам разработчиков, — это «когнитивная технология, которая может мыслить как человек» (IBM, 2017b). Watson прославился, когда стал первым искусственным интеллектом, победившим человека в игре Jeopardy!. Но на самом деле это просто пища для СМИ*. Что же де-

* Что, кстати, является еще одним примером того, как данные меняют наш способ потребления информации. Самые читаемые новостные онлайн-статьи будут вытал-

лает Watson столь привлекательным для нас? Эта технология позволяет применять данные в здравоохранении. Watson полезен прежде всего тем, что помогает врачам выявлять болезни пациентов.

Watson применяет тот же принцип, что и приложение SkinVision: собранные данные служат для диагностики — только для этого, естественно, требуются более изощренные алгоритмы. В одном удивительном случае Watson смог диагностировать редкий тип лейкемии у женщины всего за десять минут, в то время как у врачей это заняло бы несколько недель (Otake, 2016).

Все еще сомневаетесь относительно перспективы использования ИИ в медицине?

Разумеется, Watson не является решением всех наших проблем. Искусственный интеллект машин все еще может ошибаться. Но разница между машинами-врачами и людьми-медиками — это данные, и, по мере того как технология обработки растущих объемов информации совершенствуется, меняется и разница между человеком и машиной. В конце концов, люди могут поглощать информацию на конференциях, из медицинских журналов и статей, но все мы имеем ограниченную способность хранить знания. Более того, знания, которыми обладают люди-врачи, в значительной степени зависят от их жизненного опыта. В то же время врач-машина может совершенствоваться, только получая все больше данных. Благодаря мгновенному доступу к данным с других компьютеров через облако общие данные могут способствовать постановке более точных диагнозов и выполнению операций по всему миру. Благодаря экспоненциальному росту эти машины будут хранить информацию о всех видах изменений в человеческом теле, оставляя знания людей далеко позади.

Наука о данных и принадлежность

За удовлетворением потребности в безопасности (второй уровень в пирамиде Маслоу) следует потребность в принадлежности к социальной

киваться на вершину кучи, что делает это войной за самый интригующий заголовок, а не за самый убедительный контент.

среде (семья, друзья, отношения). Утверждается, что мы должны быть частью сообщества людей, которые разделяют наши интересы и видение жизни. В последние годы ощутимый разрыв между технологиями и обществом стал предметом серьезной дискуссии. Интернет часто критикуют за то, что он способствует все более изолированному существованию человека, удовлетворяя все его прихоти и потребности. Будучи любителем природы, я не стану превозносить цифровую социализацию. Тем не менее относительная доступность интернета во всем мире в любое время суток, на мой взгляд, является большим преимуществом для человеческого существования и опыта.

Более того, социальные сети, такие как Facebook, Instagram и LinkedIn, успешны не из-за удобства использования платформ, а благодаря их *данным*. Социальная сеть, на которую неохотно подписываются, вряд ли предлагает то же самое, что и сеть с большим числом подписчиков, поскольку социальная связь в конечном итоге зависит от отношений. Если нет данных, чтобы предоставить правильную информацию, будь то человеческие связи, адресованные нам изображения или новостные сюжеты по интересующим нас темам, социальная сеть окажется для нас бесполезной.

Данные позволяют сделать наш мир намного более взаимосвязанным, и это не только помогает в личных запросах, таких как поиск старых школьных друзей; они также дают возможность ученым и практикам, занимающимся схожими проблемами, найти друг друга и завязать партнерство.

Кейс: установление контактов через LinkedIn

Мне нравится использовать LinkedIn — и я думаю, что эта социальная сеть действительно научилась применять свои данные, чтобы приносить пользу как себе, так и пользователям. Быстрый переход на вкладку «Люди, которых вы знаете» — и вот уже у вас есть бесконечный список пользователей LinkedIn, с которыми вам рекомендуют установить контакт. Одни из них могут быть вашими сослуживцами, другие — бывшими однокашниками. LinkedIn использует данные, которые вы публикуете в своем профиле, — происхождение, опыт, образование, коллеги — и сопоставляет их с профилями других участников сети.

Технология LinkedIn позволила тысячам людей восстановить связи с их прошлым. И поскольку эти контакты множатся так же, как и данные сети, создается еще больше соединений. Всякий раз, когда вы подключаетесь к другому пользователю, вы выходите на связанных с ним коллег, то есть получаете соединение не только «первой степени», но и соединения «второй степени». Тем самым вы расширяете круг намного больше, чем это представляется на первый взгляд.

Для LinkedIn, как и для любых других социальных сетей, все, что необходимо, — это запрос от пользователя. Я нашел многочисленных друзей и бывших одноклассников на этом сайте, многие из которых с тех пор перешли в ту же профессиональную область, что и я. Данные соединили нас, и это открыло возможности для нового диалога старых знакомых. Осознание того, что благодаря интернету у меня сохраняются связи с друзьями и коллегами, создает чувство общности, которые не исчезают и тогда, когда мы, например, переезжаем в другой город или меняем место работы. Я нахожу эту взаимосвязь успокаивающей.

Соединяя нас с другими людьми, которые разделяют наши интересы, с которыми мы вместе учились или жили рядом, LinkedIn также может дать нам хорошее представление о работе, которая могла бы нам подойти. Когда я хотел перейти на новую работу, то начал обновлять статус на LinkedIn. Алгоритмы обработки данных этой платформы определили мои потребности в соответствии с использованными мной ключевыми словами, и именно так на меня обратили внимание рекрутеры. Еще лучше было то, что, поскольку я написал об интересующих меня предметах, алгоритмы LinkedIn подбирали мне вакансии из тех сфер, которым соответствовали мои конкретные знания. Именно так меня нашел выпускающий редактор этой книги. Как вам такая способность социальных сетей приносить счастье?

Общественное вмешательство

Хотя присутствие в онлайне может значительно улучшить как нашу личную, так и профессиональную жизнь и способствовать удовлетворению потребности в принадлежности к социальной среде, мы также должны знать о том, чем оно чревато. Одна из самых больших проблем — в том, как защитить наши данные от кражи. Кибербезопасность стала горячей темой с момента роста онлайн-банкинга, и электронная коммерция уже является *modus operandi* розничной торговли

для охвата новых клиентов. Раньше нам советовали чаще обновлять пароли, делать покупки только на проверенных сайтах, а если наши банковские реквизиты оказывались под угрозой — как можно скорее связаться с отделом банка по борьбе с мошенничеством. Учитывая, что мы все чаще осуществляем транзакции в интернете, нам стоит беспокоиться тем, как компании защищают нашу информацию.

Кейс: утечки данных и программы-вымогатели

Чем больше вы пользуетесь интернетом и чем крепче связаны с другими пользователями, тем неизбежнее увеличится объем ваших «выхлопных данных». Чем больше данных вы производите, тем более ценным источником дохода вы становитесь для компаний, продающих информацию о пользователях. Данные заменили нефть в качестве самого ценного ресурса в мире (*The Economist*, 2017).

Но когда вещи становятся ценными, они могут стать объектом кражи или злоупотребления. И, учитывая то, насколько тесно мы связаны, забота о нашей личной информации сегодня выходит далеко за рамки номеров кредитных карт. Масса личной информации размещается в интернете, и всякий раз, когда наш персональный компьютер подключен к Сети или внешнему серверу, мы рискуем, что эту информацию украдут. Чтобы увидеть потенциальный масштаб этого риска, достаточно только вспомнить глобальную кибератаку WannaCry в мае 2017 г., когда в 150 странах компьютерный червь заразил компьютеры Microsoft с целью вымогательства. Во множестве учреждений, включая FedEx* в Соединенных Штатах и министерство иностранных дел Румынии, червь WannaCry зашифровал данные пользователей — от отдельных лиц до организаций глобального масштаба, а разработчики вирусной программы требовали платы в обмен на расшифровку данных. В конечном итоге у пострадавших не было выбора, кроме как заплатить команде разработчиков за выкуп своих данных, чтобы предотвратить их уничтожение.

Такова сила данных: их кража за несколько секунд может поставить на колени целую организацию.

Еще одним примером серьезного нарушения кибербезопасности стала утечка данных Equifax. Агрегатор данных более чем 800 млн потребителей и более

* Американская компания, предоставляющая почтовые, курьерские и другие услуги логистики по всему миру. — *Прим. пер.*

88 млн предприятий во всем мире, Equifax считается одной из кредитных компаний «Большой тройки». 7 сентября 2017 г. Equifax объявила, что киберпреступники похитили идентификационную информацию компании и что эта кража могла затронуть 143 млн потребителей в США. Похищенная информация содержала имена и фамилии, даты рождения, номера полисов социального страхования, адреса и т.д. (Haselton, 2017). Учитывая, что население США в то время составляло 324 млн человек, пострадал почти каждый второй житель страны.

Рост кибербезопасности

Число и масштаб кибератак на потребителей и учреждения растут. В то же время киберпреступники становятся все более осторожными, что затрудняет даже обнаружение их местоположения. Распространение биткойна, цифровой платежной системы, позволяющей осуществлять анонимные переводы, усугубляет и без того сложную проблему поиска хакеров и привлечения их к ответственности. То, что организовать утечку информации можно из любой точки мира, не позволяет правоохранительным органам оперативно находить преступников.

Сегодня неудивительно, что специалисты по кибербезопасности пользуются высоким спросом. Такие профессионалы противостоят мошенникам и хакерам в режиме реального времени, а также проводят экспертно-криминалистический анализ после того, как произошли атаки. По мере того как меняется наше взаимодействие в интернете, как развиваются и меняются цифровые системы, люди овладевают новыми способами мошенничества онлайн и в нашем распоряжении появляются новые онлайн-средства для борьбы с ними. Специалисты по кибербезопасности должны постоянно играть в кошки-мышки, если они хотят опережать угрозы.

Что бы я посоветовал тем, кто хочет заниматься кибербезопасностью? Узнайте, как работать с неструктурированными данными, то есть с нечисловой информацией. Как правило, 80% данных компаний не структурированы (SuperDataScience, 2016). Более подробно мы рассмотрим специфику работы с неструктурированными данными в следующей главе.

Как защититься от кибератак?

Если мы используем компьютеры, подключенные к интернету или внешним серверам, и особенно — социальные каналы для обмена информацией, полностью защититься от кражи данных невозможно. Однако в наших силах более внимательно относиться к хранению и управлению данными, чтобы эффективно противостоять любой опасности. Я советую вам использовать приемы, которые применяю для защиты моих данных:

1. Храните копии всех файлов, которые вы не можете позволить себе потерять, на внешнем жестком диске или выносной памяти.
2. Регулярно копируйте жесткий диск на надежный внешний жесткий диск.
3. Присвойте ярлыки своим онлайн-аккаунтам и закрывайте все аккаунты, которыми вы больше не пользуетесь.
4. Архивируйте данные, которые вам больше не нужны, и отсоедините их от интернета. Убедитесь, что эти файлы надежно хранятся, и держите архивы в прохладном, надежном месте.
5. Не храните конфиденциальную информацию на обменных серверах облачного типа.
6. Проводите регулярные проверки программного обеспечения, чтобы обнаружить возможную утечку данных до того, как она произойдет. Вирусы-вымогатели и черви могут месяцами находиться в пользовательской системе, заражая все укромные уголки баз данных и портя даже резервные копии, прежде чем наконец зашифровать данные.

Наука о данных и признание

Потребность в признании — четвертая по важности потребность, по мнению Маслоу. Признание может быть обеспечено с помощью данных. Многие цифровые рабочие платформы помогают клиентам, агентствам и фрилансерам найти наиболее подходящего человека

для выполнения конкретной задачи, используя рекомендации и старринг — системы главных ролей. Как только проект завершен, онлайн-новые фриланс-платформы дают участникам возможность публично оценить друг друга на основе параметров, варьирующихся от доступности до качества работы. Каждая платформа имеет свою рейтинговую систему, но в целом эти данные в конечном итоге помогают клиентам найти оптимального исполнителя; также они стимулируют получающих хорошую оценку фрилансеров к продолжению работы на высоком уровне и вынуждают тех, кто получает отрицательный отзыв, совершенствовать свои профессиональные навыки. Некоторые могут быть против того, чтобы подвергнуться такой проверке, но последовательно публикуемые данные о качестве работы позволяют людям определить, в чем они преуспевают, а где им может понадобиться дальнейшее обучение.

Данные заслуживают признания

Компаниям следует подтолкнуть пользователей к тому, чтобы они включали в общие базы демографические данные о себе (такие, как возраст и местоположение). Также потребуется разработка более всеобъемлющей системы, выходящей за рамки простого метода главных ролей, и проведение по этим обзорам неструктурированного анализа, который должен дать более ценный и точный пример того, как чувствует себя пользователь. Затем данные могут быть визуализированы в облаках слов (популярные визуальные представления текста, о которых мы узнаем больше в следующей главе) или быть доступными через фильтры, применимые к демографическим данным пользователей.

Наука о данных и самореализация

Вот где начинается самое интересное (буквально). Под «самореализацией» Маслоу понимает потребность человека реализовать свой потенциал в жизни. В отличие от низших уровней иерархии, которые в значительной степени отражают врожденные потребности

всех людей, рассматриваемая здесь потребность может проявляться по-разному — ощутимо или неосознано — в зависимости от интересов человека. Потребность одного человека в самореализации может быть удовлетворена, когда он овладеет навыками рисования акварелью, а другого — когда он станет хорошим, способным убеждать своих слушателей оратором.

Кейс: игровой опыт

В конечном счете самореализация имеет отношение к потребности человека в радости. И мы уже видели, какова важность этого для индустрии развлечений. Индустрия видеоигр, в которой ворочаются миллиарды долларов, имеет очевидные связи с наукой о данных в их зависимости от технологий. Виртуальная реальность (VR) является одной из самых захватывающих областей, в которых данные специально используются для дальнейшего развития и улучшения игрового опыта. Там, где VR когда-то считалась причудой, теперь она является основным направлением в отрасли — и это в значительной степени благодаря продвинутым возможностям технологии обработки данных, например в том, что касается частоты кадров и деталей, необходимых для создания реалистичного виртуального мира. До прорыва в развитии, произошедшего в 1990-х, возможности системы автоматизированного проектирования (САПР) были ограничены отсутствием технологии его построения. Теперь данные можно использовать для создания полноразмерной виртуальной 3D-среды, в которой задействованы алгоритмы, отслеживающие ваше «местоположение» в этой среде в реальном времени, что позволяет экранам игроков подстраиваться под их взгляд с помощью 3D-очков с активным затвором и 3D-проекторов.

Именно так данные улучшают технику видеоигры. Но они также могут быть использованы для совершенствования опыта игрока путем учета того, как он ведет игру. И данные от пользователей могут быть собраны гораздо большим количеством способов, чем это возможно в других развлекательных отраслях, таких как кино. Оставляемые пользователями «выхлопные данные» охватывают взаимодействие игроков, игровое время, расходы на дополнительные игровые компоненты и активность в игровых чатах. Тем самым оптимизируются не только рекомендательные системы и реклама, но и механика игры, так как выявляются возможности сделать ее более приятной. В ход идут даже большие данные, которые создаются платформами распространения программного обеспечения и позволяют

предсказывать периоды максимальной загрузки и выбирать время, наиболее подходящее для посещения игровых серверов.

Заключительные размышления

Очевидно, что развитие науки о данных пошло на пользу огромному числу областей нашей жизни. И данные продолжают создавать пронизываемый слой между физическим и цифровым ландшафтами, переопределяя то, как мы взаимодействуем с обеими средами. Это может вызвать некоторые противоречивые мысли, но, как видно из того, как легко данные могут быть соотнесены с пирамидой потребностей Маслоу, развитие, управляемое данными, в корне облегчит человеческое существование.

Естественно, многие из этих разработок и то, как мы адаптируемся к ним, зависят от аналитика данных, поэтому в следующей главе я опишу, как можно размышлять с позиций такого специалиста. Также мы убедимся, что наше первое погружение в дисциплину должным образом направляется, и узнаем, как применить опыт, который у нас уже есть.

Мышление, необходимое для эффективного анализа данных

03

Я не утверждаю, что если вы прочтете эту книгу, то станете экспертом в области науки о данных, но, безусловно, есть способы, с помощью которых вы можете начать менять свое мышление, чтобы получить преимущество перед другими, кто тоже хочет познакомиться с этой дисциплиной. Такова цель главы 3. Всем известно, что, если вы играете на музыкальном инструменте, необходимы годы практики, прежде чем вы овладеете им на профессиональном уровне. Нужно освоить гаммы и арпеджио, ваши пальцы должны скользить по клавишам, будто они смазаны маслом, и ваши соседи, вероятно, станут протестовать против шума прежде, чем вы только осмелитесь приступить к Рахманинову. Короче говоря, чтобы превратиться в хорошего музыканта, нужны значительные инвестиции вашего времени и денег.

Наука о данных обходит стороной этот трудоемкий процесс. Даже если вы изучите только самые основные «гаммы» — например, первые несколько алгоритмов, приведенных в главе 6 «Анализ данных» (часть I), — вы все равно значительно продвинетесь на пути к работе с очень сложным материалом. И как любой, кто имеет доступ к компьютеру, также сможете познакомиться с множеством бесплатных онлайн-программ и презентаций, касающихся анализа данных (а также курсов по науке о данных). Вы почти сразу сможете начать совершенствовать вашу технику, позволив программному обеспечению выполнять за вас подготовительную часть, пока вы сосредоточены на творческой составляющей своего проекта.

Хотя я всегда призываю тех, кто планирует заниматься наукой о данных, читать и узнавать как можно больше о ней, чтобы добраться

до вершин своей профессии, я должен также подчеркнуть, что первое вхождение в предмет не должно быть ошеломляющим. Хотя и существуют некоторые предпосылки к тому, чтобы стать аналитиком данных (их мы рассмотрим более подробно в главе 10), я выбрал пять ключевых атрибутов для соответствующей настройки вашего мышления. Они позволят вам прямо сейчас приступить к освоению этой дисциплины.

1. Выберите правильное место, чтобы начать

Аналитикам данных не нужно знать всех тонкостей каждой части программного обеспечения и каждого алгоритма, чтобы разбираться в этой области. Существует огромное множество доступных программ, а алгоритмы варьируются от простейших, способных классифицировать данные, до самых сложных, использующихся в искусственном интеллекте. Когда вы в самом начале пути, то, прежде чем погрузиться в определенную область, нужно потратить время и выяснить, в какой сфере лежат ваши интересы, будь то визуализация или машинное обучение. Воздержитесь от спонтанного ответа — он не только ограничит вас на начальном этапе изучения науки о данных, но и может лишить вдохновения, если вы совершите ошибку при выборе. Многим визуализация может показаться интереснее, чем анализ, но вы должны не жалеть времени на то, чтобы понять, что требуется в каждом случае. Хорошая новость заключается в том, что к тому моменту, когда закончите читать эту книгу, вы будете гораздо яснее представлять, какая область интересует вас больше всего.

Давайте также уточним, что мы имеем в виду, говоря об ориентации на конкретную область; существует большая разница между выбором ниши, из которой вы можете совершить прыжок в своей карьере, и специализацией в ней. Последнее — опасный шаг, делать который я бы никогда не посоветовал. В конце концов, наука о данных — динамичный предмет и требует от своих практиков быть столь же динамичными в исследовании того, как решать новые проблемы в этой области. Алгоритмы меняются, программное обеспечение — тоже,

Закон Мура 2.0

Обобщим то, что мы узнали в главе 1: закон Мура является проекцией экспоненциального роста и основан на первоначальном наблюдении, что количество транзисторов в интегральной схеме будет удваиваться каждые два года. С тех пор этот закон используется для учета темпов развития (и обратно пропорциональных затрат) в области технологии и для прогнозирования того, как скоро будущие достижения могут стать реальностью. Тот факт, что каждый год у нас появляется новый iPhone с процессором примерно на 50% быстрее, чем у предыдущей модели, служит одним из таких примеров действия закона Мура.

В отличие от ситуации 30-летней давности, когда доступ к средствам обработки данных имели только сотрудники разведывательных служб и правительственных органов безопасности, сегодня даже детям дошкольного возраста доступен широкий спектр данных с лежащих в их заднем кармане ручных устройств. Закон Мура позволяет нам получить доступ к данным, исследовать и использовать их потенциал через этот взрыв технических достижений.

Одним из моих любимых примеров действия закона Мура на практике является проект «Геном человека», который был запущен в 1990 г.* Участники проекта поставили перед собой задачу определить последовательность пар оснований нуклеотидов, составляющих ДНК человека. Медленные темпы в первые годы реализации проекта вызывали обеспокоенность у тех, кто наблюдал за его развитием извне. По прошествии первых семи лет прогнозисты подвели итог — в какой части генома последовательность установлена — и предсказали, что для завершения работы потребуется еще 300 лет. Однако в этих прогнозах они не учли закон Мура. Конечно же, следующие семь лет проекта ознаменовались полным и успешным секвенированием генома — примерно на 294 года раньше запланированного срока, если принять во внимание линейную прогрессию.

* Данные из этого проекта находятся в свободном доступе по адресу www.internationalgenome.org.

и специализация в том, что в будущем перестанет существовать, не является конструктивным способом практиковать рассматриваемую дисциплину. Как мы обсуждали в главе 1, аналитики данных должны быть хорошо осведомлены о росте и переменах. Это особенно верно, если учесть, что скорость технологического развития непосредственно влияет на их работу, как это определено законом нашего старого друга Мура.

2. Напрягите творческие мышцы

Как мы узнали, массив данных будет полезен не меньше, чем аналитик данных. Для любого проекта требуется высокая степень креативности, чтобы получить максимальную отдачу от имеющихся данных. Аналитики данных должны проникнуться мышлением, позволяющим задавать правильные вопросы об интересующих их данных, и я хочу подчеркнуть здесь, что вы должны думать творчески и нестандартно — определяя далекоидущие последствия проекта через его массивы данных. В конце концов, применение данных способно дать удивительные результаты — высветить проблемы, нюансы и пробелы, о которых мы, возможно, не узнали бы без тщательного анализа данных. Это актуально для всех дисциплин и отраслей, которые используют данные для управления практикой: креативность — вклад аналитиков данных в наилучшее решение проблемы — значительно повлияет на качество выполнения задания.

Конечно, необходимый уровень творчества варьируется: для решения одних проблем достаточно традиционного подхода, а для решения других нужно что-то оригинальное. И если вы спросите меня, что лежит на дальнем конце этого спектра и что находится на переднем крае науки о данных и технологий, без тени сомнения я отвечу: искусственный интеллект.

Времена высокочувствительных роботов из «Бегущего по лезвию» придут еще не скоро, но было много ситуаций, когда роботы брали верх над людьми, играя с ними в человеческие игры.

Искусственный интеллект

С кем бы я ни разговаривал, упоминание искусственного интеллекта (ИИ) всегда вызывает интерес. Это увлекательная область развития, новости о которой обязательно попадут в заголовки. Однако ИИ полностью зависит от наличия данных и способности компьютера их обрабатывать.

Первое, о чем многие подумают при обсуждении ИИ, — это отношение к нему в голливудских фильмах, предупреждающих, что прогресс в этой области в конечном итоге приведет к нашей гибели. В «Бегущем по лезвию», экранизации научно-фантастического романа Филипа К. Дика «Мечтают ли андроиды об электрических овцах?», облик и реакции роботов («репликантов») настолько реалистичны, что в конечном итоге они становятся угрозой для существования человека. По этой причине роботов изгоняют во внеземные колонии. Однако некоторые из них возвращаются на Землю и ведут себя враждебно по отношению к нашему биологическому виду. Поскольку отличить этих роботов от людей по внешнему облику невозможно, создается машина Войта–Кампфа. Она подобна полиграфу и фиксирует ответы на ряд вопросов, специально разработанных для изучения эмоциональной реакции испытуемых. Предполагалось, что эти вопросы озадачат роботов — поскольку у них эмоции вроде бы отсутствуют — и тем самым раскроют истинную идентичность репликантов.

Реальный прототип теста известен как тест Тьюринга. Предложенный дешифровальщиком Аланом Тьюрингом в 1950-х гг. для оценки способности людей отличать машину от человека, тест оценивает ответы, полученные во время опроса. В отличие от теста Войта–Кампфа, в тесте Тьюринга два субъекта: один — робот, другой — человек, и оба они скрыты от взгляда исследователя. Последний должен определить, какой из субъектов является роботом*, — он задает обоим ряд только текстовых вопросов и оценивает, насколько их ответы похожи на те, что мог бы дать человек.

* При этом задача робота — отвечать так, чтобы его не понял исследователь. — *Прим. науч. ред.*

Кейс: Deep Blue и AlphaGo

В соревнованиях 2016 г. по игре в го (очень популярная в Восточной Азии абстрактная стратегическая настольная игра, в которой участвуют двое) машине, известной как AlphaGo и созданной дочерней компанией Google DeepMind, удалось победить 18-кратного чемпиона мира Ли Седоля в четырех из пяти игр.

Вы можете не считать это каким-то грандиозным достижением, вспомнив знаменитую шахматную партию, сыгранную русским гроссмейстером Гарри Каспаровым и Deep Blue, компьютером, специально разработанным IBM. Deep Blue выиграл, и это случилось еще в 1997 г. Но даже несмотря на то, что робот добился успеха почти за 20 лет до успеха AlphaGo, результат последней представляет для нас особый интерес.

Игра в шахматы полностью основана на логике. Цель Deep Blue состояла в том, чтобы безупречно соблюдать эту логику и ждать, пока противник допустит ошибку. Люди совершают ошибки, машины — нет. В отличие от шахмат игра в го основана на интуиции. По сравнению с логикой, которой руководствуется компьютер, интуиция — гораздо более сложный феномен: она требует, чтобы машина развивала внутренние знания об игре, которые не могут быть просто запрограммированы в ней*.

В го игроки перемещают черные и белые фишки по доске с разметкой 19×19 клеток. Цель игры — захватить большую площадь, чем противник. AlphaGo первоначально получила обширную базу данных — около 30 млн сделанных людьми ходов, проанализированных с помощью комбинации машинных алгоритмов и методов свободного поиска. После того как значительное количество игр было сыграно против соперников-людей и собрано достаточно знаний о поведении противников, AlphaGo миллионы раз сыграла сама против себя, чтобы еще больше улучшить результаты. (Это тип обучения с подкреплением, о котором я расскажу более подробно в главе 6.) Только после того, как этот период обучения завершился, создатели машины выставили ее против лучших игроков мира. От шахмат до го искусственный интеллект прошел значительный путь,

* Генеральный директор подразделения Google DeepMind Демис Хассабис определяет интуицию как неявное знание, которое приобретается через опыт и не является сознательно выраженным или даже доступным, поэтому мы не можем получить доступ к этому знанию сами и, конечно, не можем передать его другим.

обучаясь через действия и наблюдения, а не только применяя математическую логику*.

В этот момент вы можете подумать: «Победа ИИ в шахматах и го впечатляет, но как все это относится к бизнесу?»

Применение искусственного интеллекта не ограничивается победами над людьми в игре го. Та же компания DeepMind разработала искусственный интеллект, чтобы помочь Google лучше управлять охлаждением в их обширных центрах обработки данных. Система смогла последовательно достигнуть поразительного 40%-ного сокращения количества энергии, используемой для охлаждения. Это не только создает огромный потенциал для экономии в компании, но также означает повышение энергоэффективности, сокращение выбросов и в конечном счете — вклад в решение проблемы изменения климата (DeepMind, 2016). Если это не творческий подход к решению проблем бизнеса, то я не знаю, что им является.

3. Используйте свое прошлое

Как я уже говорил в пункте 1, истинная красота науки о данных заключается в том, что в отличие от многих других дисциплин для ее освоения не потребуются годы практики. Читатели, которые только начинают заниматься наукой о данных, не должны чувствовать себя в невыгодном положении относительно сверстников, которые, возможно, работали с данными и изучали их всю жизнь. Опять же, все, что вам нужно, — это небольшое изменение в мышлении — сосредоточьтесь на том, что вы *знаете*, а не на том, чего не знаете. Используйте и свои углубленные знания другого предмета, и любые навыки, которые вы, вероятно, получили как профессионал и/или студент.

* В октябре 2017 г. Google DeepMind анонсировала AlphaGo Zero. Его особенность заключается в том, что он вообще не использует никаких человеческих данных, а скорее полностью учится на собственной игре (DeepMind, 2017). Эта новая версия настолько мощная, что победила первую Alpha Go в 100 играх. Как это коррелирует с экспоненциальным прогрессом в науке о данных?

Глубокие знания

Мало того что в науку о данных несложно вникнуть — занявшись ею после освоения какой-либо другой дисциплины, вы получаете *преимущество*. Вот где творческий стержень науки о данных может проявить себя еще раз. Возьмем в качестве примера писателей-профессионалов. Если писатель потратил все свои усилия только на изучение того, как и что писать, и у него не было времени на расширение своего кругозора, на прочтение множества книг по самым разным вопросам, то у такого писателя не хватит знаний и опыта, чтобы опираться на них в работе. То же самое верно для науки о данных: те, кто изучал только ее всю свою жизнь и имеет ограниченный профессиональный или личный опыт в других сферах, будут подходить к любому проекту однобоко.

Итак, предположим, что лингвист решил заняться наукой о данных. Он будет иметь значительное преимущество перед другими аналитиками данных в связанных с лингвистикой проектах. Это правда: назовите любую профессию, и я расскажу вам, как применить в ней науку о данных. Аналитик данных с опытом в лингвистике, например, мог бы выиграть от получения доступа к материалам из Международного архива диалектов английского языка, в котором хранятся голоса тысяч участников со всего мира, и использовать эти звуковые файлы для составления диалектной карты мира. «Сырой» аналитик данных может поэкспериментировать с материалом, но специалист по данным с *правильным прошлым* задаст правильные вопросы, чтобы получить действительно интересные результаты. Скажем, Вест-Индия, известная лингвистам распространенным там необычным сленгом, может быть взята в качестве объекта первоначального исследования, результаты которого заложат основы для дальнейшего изучения поколенческих, этнических и гендерных различий в речи.

Стать специалистом в области науки о данных не означает разворот на 180° по отношению к тому, что вы узнали и освоили раньше. Как раз наоборот. Иногда самые интересные для вас проекты будут находиться «рядом с домом». Подумайте о проблемах, с которыми вы сталкиваетесь на своем рабочем месте: есть ли способ решить их с помощью данных?

Гибкие навыки

Хотя это, несомненно, полезно, вы необязательно должны быть экспертом в какой-то области, чтобы иметь фору в науке о данных. Даже гибкие, широко использующиеся навыки, такие как работа в команде и опыт публичных выступлений, могут значительно помочь вам. Они принесут даже больше пользы, чем глубокие знания, тем, кто, еще не успев получить достаточный жизненный опыт или образование. Подумайте о своих навыках: вы легко общаетесь? Можете ли вы адаптировать устоявшиеся решения к различным ситуациям? У вас эстетический вкус? Вы нестандартно мыслите?

Я пришел в науку о данных, будучи специалистом в области финансов, но, хотя мои знания, несомненно,годились мультинациональной консалтинговой фирме Deloitte, думаю, что в конечном итоге мне помогли гибкие навыки, которые я приобрел гораздо раньше, еще в школьные годы. Кроме того, начиная заниматься наукой о данных, я хорошо понимал, как визуализировать результаты проектов эстетически привлекательным образом. В детстве я жил в Зимбабве, где дважды в неделю изучал изобразительное искусство. Я приобрел только базовые навыки в рисовании и научился лепить забавную глиняную посуду, но, хотя курс, возможно, и не сделал меня преемником Жоана Миро*, он научил меня тому, как цвет, эстетика и положительные психологические эффекты могут повлиять на мой итоговый рабочий отчет.

После того как несколько лет спустя я вернулся в Россию, мне преподавали — в трех разных школах — совсем другое, в основном точные науки. Это научило меня академической строгости, которая пригодилась в грядущие годы в университете, но привело к нехватке необходимых социальных навыков. Будучи почти неисправимым интровертом, я работал над собой, чтобы приобрести некоторую уверенность в себе и развить способность к общению — качества, которые, как я знал, мне понадобятся. Я нашел книгу по самопомощи, в которой было все, что мне требовалось знать о том, как выбраться

* Каталонский художник-абстракционист, прославившийся в том числе керамическими работами. — Прим. пер.

из своей раковины. Упражнения в ней были немного необычными (общаться, лежа посреди оживленной кофейни, или завести непринужденный разговор с людьми в общественном транспорте), но для меня они сработали. Эти усилия, первоначально, возможно, мотивированные юношеским стремлением к участию в университетских мероприятиях и спортивных командах, позже помогли мне зарекомендовать себя коммуникативным человеком, что оказалось привлекательным на моей работе, где были нужны аналитики данных для предоставления отчетов широкому кругу лиц, имеющих отношение к деятельности всей компании.

Это еще один важный фактор для аналитиков данных: если вы хотите получить возможность запустить проект по обработке и анализу данных, вам придется научиться разговаривать с нужными людьми. Это часто будет означать расспросы *вне* вашей команды и потенциальной зоны комфорта. Данные ничего не скажут вам, если вы не зададите правильные вопросы, поэтому ваша работа — выйти «в свет» и получить ответы от людей, которые внесли свой вклад в вашу базу данных.

В обоих случаях, которые мы здесь обсуждали, — используете ли вы связанное с углубленными знаниями преимущество для поиска информации и применяете ли гибкие навыки для получения ответов от людей, — вы, вероятно, сталкиваетесь с данными, которые не являются числовыми и истинность которых поэтому зависит от контекста и субъективности анализа. Информация такого рода — мы называем ее неструктурированными данными — может быть письменным ответом, либо записанным видео- или аудиоинтервью, либо изображением. По причине того, что неструктурированные данные нельзя оценить количественно, компании часто предпочитают приглашать для их анализа специалистов в соответствующих областях.

4. Практика ведет к совершенству

Одним из замечательных аспектов науки о данных является то, что существует множество бесплатных материалов с открытым исходным кодом, которые позволяют легко продолжать практиковаться.

Аналитика неструктурированных данных

Неструктурированная аналитика работает, как вы уже догадались, с неструктурированными данными, которые составляют большую часть информации в мире. Давая определение неструктурированным данным, проще сказать, что это все, что не относится к структурированным данным (числовой информации). Это может быть текст, аудио, видео или изображения. Название объясняется тем, что этот вид данных нельзя непосредственно преобразовать в массив данных — их необходимо сначала подготовить, а поскольку неструктурированные данные зачастую нельзя автоматически перевести в исчисляемые, то в их анализе неизбежна некоторая степень субъективности. В связи с этим неструктурированная аналитика крайне важна для любого исследователя данных.

Классическим примером неструктурированной аналитики является работа с качественными опросами, которые дают данные в текстовом или ином нечисловом формате. В прошлом эти данные должны были быть преобразованы в числовую форму, прежде чем их можно было понять с помощью аналитических инструментов. Это означало, что любые вопросы, которые не предполагали множественного выбора или одного ответа — и поэтому не могли быть легко перенесены в числовой формат, — требовали от аналитика данных вручную производить численную классификацию каждого ответа.

Например, на вопрос о том, чем наслаждался посетитель Йеллоустонского национального парка во время своего пребывания в нем, можно было получить ряд ответов, включая «полевые цветы», «пикники», «занятия живописью», «наблюдение за птицами», «греблю на каяке», «отличный отель с завтраком» и т.д. Аналитик данных должен был бы прочитать все эти результаты, а затем вручную сгруппировать их в категории, которые, по его мнению, были значимыми, такие как «природа», «деятельность», «экскурсии» и «отдых». Не всегда легко сгруппировать ответы по категориям, так как здесь не исключен субъективный подход.

Вы можете себе представить, что перевод этих ответов в числа в лучшем случае представлял итоговый массив данных в немного искаженном виде.

Сегодня методы сортировки результатов по контексту кардинально изменили то, как мы проводим исследования, и новые алгоритмы в этой области помогают нам точно работать в том числе и с изображениями. Аналитики данных признают наличие проблем в методах организации качественных данных и прилагают согласованные усилия для обработки значений, которые нелегко преобразовать в цифры. Полученные алгоритмы позволяют давать гораздо более точные прогнозы, чем было возможно ранее. Теперь мы можем рассматривать слова аналогично числовым данным, например обучая аналитические инструменты идентификации вспомогательных глаголов, а также идиоматических выражений, которые имеют отдаленное отношение к заданному ключевому слову. Это позволяет машине исследовать текстовые данные куда более качественно. Здесь может прийти на ум анализ литературных произведений с помощью цифровых гуманитарных наук, но это лишь мизерная доля того, что могут делать машинные алгоритмы в этой области. Применения неструктурированной аналитики выходят далеко за рамки академической сферы и простираются в мир коммерции. Даже в криминалистике машины теперь могут просматривать письменные сообщения подозреваемых с целью установить особенности поведения, которые детектив мог не заметить.

Вы можете подумать, что люди всегда будут действовать эффективнее машин при изучении средств массовой информации: большинство из нас все еще считает, что мы всегда будем лучше понимать более широкую контекстуальную среду. Как компьютер может распознать период искусства, или стаю чаек, или эмоции лучше, чем человек? На самом деле машины уже давно могут давать ошеломляюще точные прогнозы в отношении нечисловых данных. Еще в 2011 г. исследование, проведенное Институтом нейроинформатики Рурского университета в Бохуме и кафедрой компьютерных наук

Копенгагенского университета, показало, что машины могут превосходить людей в выполнении даже таких сложных задач, как идентификация дорожных знаков (Stallkamp et al., 2012). Для этого исследования команда показала испытуемым машинам и людям фотографию, разделенную на квадраты. Задача состояла в том, чтобы определить, на каких квадратах (если таковые имелись) есть полное или частичное изображение дорожного знака. Возможно, вы видели эти тесты в интернете — в настоящее время они используются для дополнительной проверки безопасности перед входом пользователя на сайт и специально разработаны, чтобы лишить роботов доступа к защищенным данным. Результаты этого исследования показывают, что мы уже не в состоянии предотвратить захват данных искусственным интеллектом.

Облака слов

Я вижу, что облака слов часто используются в публичных презентациях, и подозреваю, что причина в том, что они искусно и содержательно сочетают изображение с текстом. Облака слов (или облака тегов) — это популярные способы визуализации текстовой информации, и если вы еще не используете их в своих презентациях, то захотите, узнав, как они работают. Создатель облака слов берет набор наиболее часто используемых слов из фрагмента анализируемого текста и группирует их в одном изображении, обозначая порядок их важности размером шрифта, а иногда также и цветом.

Облака слов, естественно, можно использовать для выделения терминов, которые чаще всего встречаются в тексте, будь то пресс-релиз или литературное произведение. Они также могут быть применены к данным опросов, что делает их очень простым, но эффективным способом показать пользователям ключевые понятия или ощущения, связанные с заданным вопросом. Таким образом, их эффективность связана с многофункциональностью и определением ключевых или наиболее значимых слов во всем, что содержит текст: метаданных, романах, докладах, анкетах, эссе или исторических записях.

В интернете есть много простых генераторов облака слов, где вы можете поиграть со шрифтами, макетами и цветовыми схемами. (Облако слов на рис. 3.1, например, было сгенерировано на основе анализа текста введения этой книги с использованием www.wordclouds.com.) Они гораздо более привлекательно выглядят, чем упорядоченные списки. Обратитесь к ним при подготовке своей следующей презентации; вы удивитесь тому, как легко окажется запустить дискуссию (подробнее о средствах визуальной аналитики см. главу 6).



Рис. 3.1. Облако тегов, созданное на основе введения

Наука о данных значительно облегчила компаниям доступ к средствам массовой информации и их анализ. Большинство владельцев бизнеса и маркетологов знакомы с SurveyMonkey — онлайн-провайдером бесплатных анкет, который обрабатывает сведения, полученные из опросов, с помощью своих инструментов анализа данных. Пользователи получают доступ к потребительским данным в режиме реального времени, а ответы из анкет участников визуализируются в виде простой графики и пользовательского дашборда. На момент написания этой книги компания может предоставлять результаты в режиме реального времени, составлять пользовательские отчеты в виде диаграмм и графиков, осуществлять фильтрацию данных, выявляя демографические тенденции, а также проводить текстовый анализ, давая пользователям наиболее релевантные текстовые данные из опросов в виде облака слов.

Новички в какой-либо дисциплине склонны месяц за месяцем изучать теорию вместо того, чтобы настроить свое мышление на применение полученных знаний на практике. В качестве упражнения просто введите в поисковике слова «бесплатные массивы данных»/«free datasets» — и найдете множество сайтов, которые позволяют скачать их CSV-файлы (файлы для хранения табличных данных), готовые для анализа. Учитывая огромное количество и диапазон данных, от космических исследований NASA до комментариев Reddit или даже спортивных данных (баскетбол, футбол, бейсбол), я уверен, что вы найдете что-то ценное и интересное*.

В то время как лучшие инструменты анализа в настоящее время небесплатны для пользователей, все большее количество программного обеспечения либо имеет открытый исходный код, либо находится в свободном доступе в интернете. Если бы вы были художником, это походило бы на бесконечный запас мольбертов, красок и холстов.

Я настоятельно призываю вас использовать эти общедоступные массивы данных для проверки своих навыков и проведения собственных анализов. В практике нет кратчайшего пути. Многое из того, что вы делаете, особенно на начальном этапе, будет включать в себя пробы и ошибки. Лучший способ приучить себя отстраненно думать о решении проблем с помощью данных — повысить свою открытость различным сценариям, другими словами, различным массивам данных.

С чего начать? Лучший выбор может быть прямо под носом. Я ожидаю, что многие читатели окажутся владельцами бизнеса или сотрудниками компании, которая рассчитывает использовать данные в ближайшем будущем. Те из вас, кто так или иначе работал с какой-то компанией, в какой-то момент столкнутся с бизнес-аналитикой.

Бизнес-аналитика vs наука о данных

Если вы уже использовали бизнес-аналитику (БА) на своем рабочем месте, значит, вы уже кое-что умеете. С помощью БА вы должны определить бизнес-вопрос, найти соответствующие данные, визуализировать

* Мы также предлагаем множество бесплатных массивов данных для наших студентов на www.superdatascience.com.

и представить их убедительным образом инвесторам и заинтересованным сторонам. Это уже четыре из пяти этапов процесса изучения данных, к которым мы вернемся во второй и третьей частях. Основным исключением является то, что БА не проводит детального, исследовательского анализа данных. Она просто *описывает* то, что произошло, в процессе, который мы называем «описательная аналитика».

Наука о данных дает нам основу для ответа на дополнительные вопросы, связанные с массивом данных компании, а также для прогнозирования и идей по улучшению. У технологической исследовательской фирмы Gartner есть модель для разделения науки о данных на четыре типа, и, если бизнес-аналитика соответствует первому типу анализа, наука о данных может помочь поставить галочки для трех остальных (рис 3.2).

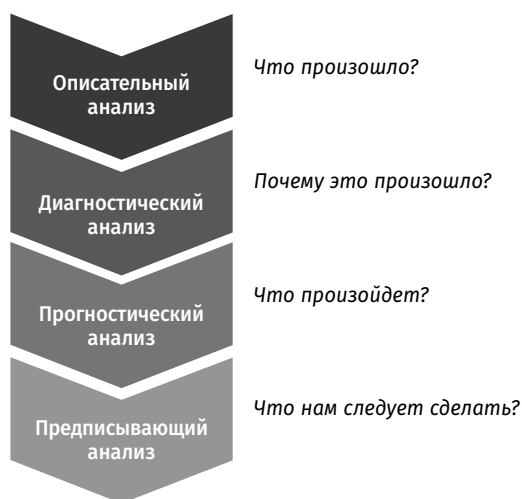


Рис. 3.2. Шкала аналитических значений

Это хорошая новость, но, если вы будете заниматься бизнес-аналитикой без учета принципов науки о данных, это может в конечном итоге помешать вашему прогрессу. Вы будете использовать данные для составления первого бизнес-отчета — но ведь владельцам бизнеса часто нужны отчеты на регулярной основе. В таком случае данные обычно отходят на второй план: все внимание приковано к конечным результатам.

Это одна из проблем БА — данные часто оказываются вторичны по отношению к содержанию обновленного отчета. Но данные *должны* быть в центре любых результатов и идей, которые связаны с бизнесом, — для каждого отчета, который мы составляем, нужно заранее провести анализ данных — иначе мы ограничимся изучением лишь тех из них, что присутствовали в предыдущем исследовании.

Цепляться за БА может быть заманчиво, когда вы или ваша компания работали таким образом в течение многих лет, но наука о данных предлагает гораздо более впечатляющий набор инструментов — образно и буквально — для анализа. Благодаря ей разрабатываются и применяются различные аналитические программы и формируется процветающее онлайн-сообщество аналитиков данных,

Все, что, как вам кажется, вы знаете, — неверно

Все мы рано или поздно сталкиваемся с Excel. Она стала одной из самых важных программ для корпораций, и большинство таблиц существуют в формате XLSX. Тем не менее для Excel характерна тенденция чрезмерного упрощения, и поэтому у вас может сложиться искаженное впечатление о данных. Если вам знакомо только представление данных в Excel, вы должны быть готовы изменить свое восприятие аналитики.

Мы подробно рассмотрим трудности с Excel в главе 5 «Подготовка данных», а здесь лишь отметим: в программном обеспечении, возможно, нет типов данных. Мы, конечно, не имеем дело с ними напрямую, а это означает, что в электронной таблице неподготовленного человека строки, формулы и визуальные эффекты окажутся перепутаны. Несмотря на то что Excel выглядит как таблица, мы можем вставлять числа, слова, ссылки и дроби в одни и те же колонки, тем самым смешивая все типы данных без разбора. Ни один инструмент науки о данных не позволит вам смешивать данные и логику — проблема, которую мы рассмотрим в главе 5. В любой системе управления базами данных логика и данные должны рассматриваться отдельно.

Будьте готовы использовать программу, которая не является Excel. На мой взгляд, одни из лучших программ для анализа массивов данных — R и Python.

работающих с открытыми исходными кодами для того, чтобы усовершенствовать процесс и поделиться своими достижениями. Возможность использования этих инструментов избавляет человека от необходимости искать информацию вручную, позволяя сосредоточиться на преодолении узких мест, раскрытии возможностей продаж и оценке работоспособности бизнес-подразделения. К сожалению, традиционная зависимость БА от Excel может научить вас плохим привычкам.

5. Помните об этике

Морозным февральским утром, задолго до того, как проснулся любой здравомыслящий человек, мне позвонили из полиции Квинсленда. Все еще сонный и едва ворочая языком, я пробормотал: «Да, я Кирилл Еременко; да, я нахожусь у себя дома в Брисбене; да, у моего байка тот номерной знак, который мне сейчас зачитали. Так в чем проблема?» Меня спросили, пользовался ли кто-нибудь, кроме меня, моим мотоциклом и знаю ли я, где он находится. Последний вопрос вернул меня в сознание и заставил слететь с лестницы в гараж.

С облегчением я убедился, что предмет моей гордости и радости все еще там. Но оставался вопрос: если все, о чем они меня спрашивали (а они спрашивали и обо мне), находилось на своем месте, то какое дело было полиции до всех этих подробностей?

Они сказали, что заметили мотоцикл с моим номерным знаком, скрывающийся от полиции в Голд-Косте, пляжном городке недалеко от Брисбена. Учитывая, что мой мотоцикл был на месте, они предположили, что номерной знак, должно быть, подделали, — и позже обнаружилось, что так оно и было.

Представьте на мгновение, что мой байк действительно украли. Как бы я смог доказать, что это не я скрывался от сотрудников правоохранительных органов? В ту ночь я был один, и про алиби не могло быть и речи. С точки зрения полиции, это, безусловно, мог быть я, особенно принимая во внимание, насколько трудно подделать номерной знак в такой стране, как Австралия, где подобные вещи жестко контролируются.

Даже несмотря на то, что в начале разговора я не знал, был ли мой мотоцикл украден, я понял, что меня совсем, даже на секунду, не беспокоил вопрос алиби во время этого телефонного допроса, потому что я знал, что не сделал ничего плохого. Я не сомневался, что технологии послужат мне как свидетели. Большую часть времени я ношу с собой телефон, заряжаю его рядом с кроватью, и любые действия, которые я выполняю с ним, регистрируются. Это напомнило мне, как в Deloitte я работал в отделе финансовых расследований (форензик). Мы разбирали бесчисленные ситуации, когда люди утверждали, будто они делали что-то или находились в определенном месте,

Этическая цена данных

Мы знаем, что данные могут причинить вред, о чем свидетельствуют бум конференций и учреждений, занимающихся изучением последствий технологического развития для этики и кодексов поведения человека. Кто имеет доступ к нашим данным? Должен ли вообще существовать доступ к ним?* Как мы видели, данные открывают перед нами новые способы работы, жизни, исследований, ведения войны — и делают это с невероятной скоростью.

Возьмем 3D-печать. По мере снижения стоимости разработки таких принтеров число людей, имеющих доступ к новой технологии, будет увеличиваться. Коммерческие 3D-принтеры в настоящее время производят игрушки и игры, но они также могут печатать любое количество потенциально опасных предметов — нужна только модель данных. Одного этого, безусловно, достаточно, чтобы вызвать обеспокоенность, особенно с учетом непропорционально высоких темпов технологического развития и нашей неспособности принимать законы и обеспечивать защиту от негативных последствий. Сможем ли мы когда-нибудь надеяться на то, что будем поспевать за таким быстрым темпом перемен?

* Дополнительные сведения об управлении данными см. в докладе, представленном Британской академией и Королевским обществом (2017), а также в серии показательных выступлений в Британской академии в рамках Сезона робототехники, ИИ и общества (British Academy, 2017, запись доступна в интернете).

но их телефоны рассказывали совсем другую историю. Эти записи использовались в качестве доказательств, потому что данные, полученные благодаря мобильным устройствам, камерам видеонаблюдения и т.п., не лгут.

Дело в том, что данные *могут* помочь. Они могут служить вашим алиби. Они могут выступать в качестве доказательства по уголовным делам. Многие считают, что данные могут только навредить, — но вы не слишком далеко продвинетесь в нашей дисциплине, если будете думать о себе как о злодее. Небольшое изменение в том, как вы рассматриваете науку о данных и ее функции, побудит вас искать новые способы совершенствования своей трудовой деятельности с помощью данных, вместо того чтобы чувствовать, что нужно доказывать свою профессиональную состоятельность коллегам.

Злонамеренное и неправильное использование данных

Один из самых острых вопросов в дискуссии вокруг технологий и этики связан с границами доступа машин к информации (Mulgan, 2016). По мере того как возможности роботов в обработке данных увеличиваются, машины скоро будут способны регулировать информацию способом, существенно превышающим возможности человека.

Информация всех видов становится оцифрованной. Хранение ее в цифровом, а не физическом формате превращается в норму. Исторические артефакты оцифрованы, книги и журналы доступны в интернете, а личные фотографии загружаются в социальные облака. В конце концов, информация намного сохраннее, когда находится в электронном виде: она не боится времени, ее можно копировать, а контент — выложить для общего пользования и установить связи между соответствующими элементами. Конечно, цифровые данные не полностью защищены от повреждений. Они могут пострадать или потеряться, но в итоге менее подвержены порче, чем данные, которые хранятся только в материальном виде.

Тот факт, что в интернете так много информации — как по охвату, так и по глубине, увеличивает потенциал машин, которые имеют доступ к этим данным, и расширяет разрыв между возможностями человека и компьютера.

Компьютеры не достигли пределов своих возможностей в обработке данных — но мы достигли. Машины ждут только трех вещей: доступа к данным, доступа к более быстрому оборудованию и доступа к более продвинутым алгоритмам.

Когда эти три условия будут соблюдены, польза и вред от машин, которые могут регулировать количество доступных им данных, станут только делом времени. И это уже закладывает основы для мощного оружия, будь то анализ поведения в интернете или маскировка под человека на сайтах социальных сетей в целях пропаганды. Если верить футурологу Рэймонду Курцвейлу, предсказавшему, что к 2029 г. компьютер пройдет тест Тьюринга, то предоставление машинам неограниченного доступа в интернет может сделать доступ к данным самым мощным инструментом манипуляций.

Мы должны также понимать, что заботы одного поколения обязательно станут заботами другого. Если мы беспокоимся о том, как информация о нас собирается, хранится и используется, то, вероятно, это не будет иметь значения для молодого поколения, выросшего

Почему бы нам просто не остановить время?

Возвращаясь домой после вечера, проведенного в центре Брисбена, я невольно оказался втянут в горячий разговор с таксистом. Он, по-видимому, негативно воспринял информацию о том, что я работаю аналитиком данных, и обвиняющим тоном заговорил о неблагоприятных для будущего последствиях моей деятельности. Опасаясь худшего, таксист жестом указал на ночное небо и спросил меня или небеса: «Почему бы просто не остановиться там, где мы находимся, прямо сейчас?»

Это просто невозможно. В нашей природе заложено стремление исследовать мир и продолжать расширять свои горизонты. Для взволнованного таксиста было естественно переживать по поводу того, как данные и алгоритмы их обработки станут использоваться в перспективе. Но тревога о том, что может произойти или не произойти, будет только сдерживать нас — пагубный сценарий, особенно с учетом того, что, пока мы паникуем, технологии продолжают развиваться.

с этой технологией. Изменение нашего взгляда на то, что мы считаем нормой, отражается в нашем подходе к сбору и обработке данных. Рассмотрим случай хранения cookie-файлов в интернете. Многие сайты предпочитают собирать данные от пользователей. Эти данные называются файлами cookie. Информация записывается в файл, который хранится на компьютере пользователя и открывается при каждом следующем посещении сайта. Файл cookie может содержать имя пользователя, адреса посещенных сайтов и даже рекламу сторонних ресурсов — все это помогает сайту адаптироваться к потребностям посетителей.

Кейс: файлы cookie в интернете

Вам может показаться знакомым следующее заявление: «Чтобы этот сайт работал должным образом, мы иногда размещаем небольшие файлы данных, называемые cookie, на вашем устройстве. Большинство крупных сайтов поступают так же». Это уведомление Европейской комиссии (ЕС), которая постановила, чтобы каждый европейский сайт, использующий файлы cookie, сообщал посредством всплывающего окна или иным образом, что он записывает данные пользователя. Те, кто желает продолжать пользоваться сайтом, могут либо сразу согласиться, либо узнать больше, прежде чем принять эти условия*. Закон был принят в то время, когда люди были обеспокоены тем, что их конфиденциальность нарушается компаниями, использующими файлы cookie для отслеживания просмотренных страниц, взаимодействий и многого другого.

С тех пор тревоги, связанные с этическим аспектом использования cookie, медленно, но верно улеглись. Никого больше не пугают файлы cookie, и уж точно — не миллениалов: мы привыкли к этим файлам как к неотъемлемой части нашей онлайн-жизни. Другими словами, озабоченность по поводу файлов cookie снизилась, и поэтому требование, чтобы на сайтах компаний содержалось четкое предупреждение о сборе данных, касающихся пользователей, будет постепенно отменяться с начала 2018 г.**

* Возможны исключения. Руководство поставщика информации о том, как подготовить согласие пользователя на веб-сайтах, доступно на сайте Европейской комиссии: <http://ec.europa.eu>.

** Пока что сайты уведомляют о том, что используют файлы cookie. Нельзя сказать, что законодательство в области сбора и хранения данных либерализуется, — наоборот, в ЕС был принят Общий регламент по защите данных (GDPR), обязывающий интернет-

Cookie — это один из примеров того, как сбор данных становится частью нашего общества. То, как большинство миллениалов используют социальные сети — например, свободно выражая свое мнение, общаясь в чате, загружая свои фотографии, отмечая друзей, — должно показать, что их мир обособлен от мира беби-бумеров и они иначе (как правило) ведут себя в интернете. Я не считаю этические соображения просто неудобными препятствиями, которые аналитик данных может предпочесть игнорировать. Но я задаю вопрос читателю: действительно ли мы должны подавлять развитие технологий, исходя из наших сегоднешних опасений? Или же нам следует стремиться к установлению баланса между темпами технологического роста и темпами разработки соответствующих этических принципов?*

Подготовьтесь к изучению второй части

Будем надеяться, что вы уже нашли что-то в своем личном и/или профессиональном опыте, что можно применить в вашей работе с данными. Отметьте навыки, которые вы можете использовать, напишите их в черновике резюме — работодатели ищут аналитиков данных, и вам существенно помогут свидетельства того, что ваше мышление изменилось и стало таким, какое необходимо для профессионалов в области данных.

ресурсы в подробностях сообщать, какую информацию они собирают и хранят. — *Прим. науч. ред.*

* Если вам кажется, что ваш проект в области науки о данных не вполне отвечает этическим нормам, я бы предложил найти или разработать этические рамки, которых ваша компания может придерживаться. Могу особенно рекомендовать документ «Этические принципы использования данных» (Data Science Ethical Framework) правительства Великобритании (UK Cabinet Office, 2016), который доступен в интернете.

ЧАСТЬ ВТОРАЯ

«Когда и где я могу получить их?»

Сбор и анализ данных

Практически в любой сфере жизни нас часто больше всего возбуждают самые сложные задачи. И проекты, в основе которых лежит использование данных, порой ставят перед нами именно такие цели. Нужно задать *новые* вопросы данным, так как от аналитиков данных всегда ждут решения *проблемы*. Когда я начинаю новый проект, мне нравится думать, что я веду разговор с данными; я общаюсь с ними, чтобы быть уверенным в том, что смогу представить их в полном и достоверном виде клиенту или участникам проекта. По моему собственному опыту и опыту моих коллег, окончательные результаты часто открывают глаза, приводя к значительным изменениям во всех учреждениях, — от тех, что занимаются практической работой, до организационных структур. Некоторые из этих результатов могут быть непосредственно связаны с бизнес-проблемой, которую вам было предложено решить, а другие способны осветить такие аспекты деловой активности, к которым организация прежде не имела доступа.

Значит, у данных есть потенциал. Это делает их столь захватывающими. Они всегда сообщают нам *что-то*, будь эта информация новой или нет. Они дают шанс продолжать изучать возможности и тем самым получать различные результаты — а для этого надо задавать различные вопросы о данных, преобразовывать их с помощью различных методов и применять к ним различные алгоритмы.

Процесс анализа и обработки данных



Из-за огромного потенциала данных доступ к ним может быть затруднен, особенно если это большой массив, который содержит различные виды данных, или если компания, для которой вы работаете, просто не знает, какие данные у них собраны. Именно здесь требуется анализ данных. Он предлагает надежную и здравую технологию для любого типа проекта, связанного с данными, независимо от объема и вида доступных сведений, и призван помочь вам выстроить свой

проект от его концепции до формы представления заказчику. Первый разработанный Джо Блишштайном и Ганспетером Пфистером процесс анализа данных ведет нас через каждый этап проекта, с момента, когда мы впервые размышляем, как подойти к данным, до оформления результатов ясным и эффективным образом.

Процесс состоит из пяти этапов:

1. Сформулируйте вопрос.
2. Подготовьте данные.
3. Проанализируйте данные.
4. Визуализируйте выводы.
5. Представьте выводы.

Каждый из этапов добавляет к вашему массиву данных то, что мне нравится называть «слой интереса». Хотя к некоторым из этих этапов можно возвращаться в ходе процесса, прохождение их в линейном порядке уменьшит вероятность ошибки на более позднем этапе проекта и поможет определить, на каком шаге произошел сбой.

Поскольку этот процесс является неотъемлемой частью каждого проекта в области науки о данных и поскольку каждый этап требует различных навыков, мы будем рассматривать этапы отдельно во второй и третьей частях книги. Вторая часть посвящена первым трем этапам. Эти первые три шага позволят нам: 1) сформулировать обоснованный вопрос или серию вопросов, на которые необходимо ответить с помощью данных; 2) собрать массив данных таким образом, чтобы он отвечал на поставленные вопросы, и 3) получить ответ из массива данных путем анализа или прогнозирования. На мой взгляд, эти этапы потребуют от вас наибольшего вклада. Если вы проделаете всю предварительную работу, то визуализировать и представить выводы будет просто, потому что вы уже достигнете целей вашего проекта.

Аналитик данных, частный детектив

Сегодня в нашем распоряжении невероятное количество данных. Подумайте о количестве комбинаций, которые можно получить

с помощью колоды из 52 игральных карт. Просто перетасуйте колоду — крайне маловероятно, чтобы кто-то еще на протяжении человеческой истории получил такой же порядок карт. Начало работы с данными похоже на то, как если бы вам вручили колоду игральных карт, — возможностей для вариаций, с которыми можно работать, иногда больше, а иногда меньше, но их всегда множество. Как только вы установили некоторые основные правила (для карт это означает игру, для науки о данных — гипотезу и алгоритм), вы действительно можете начинать работу. Определение вопроса помогает построить и спланировать подход к данным, гарантирующий, что мы получим наиболее релевантные результаты.

В «Скандале в Богемии» Шерлок Холмс говорит доктору Ватсону: «Теоретизировать, не имея данных, опасно. Незаметно для себя человек начинает подтасовывать факты, чтобы подогнать их к своей теории, вместо того чтобы подтвердить факты теорией». Холмс предостерегает Ватсона от того, чтобы строить догадки в отсутствие подтверждающих их правильность доказательств. Но то, что Конан Дойл также подчеркнул здесь, — это необходимость сделать шаг назад, прежде чем погрузиться в проблему и сформулировать какие-либо предположения или найти решение. Имея дело с данными, мы располагаем преимуществом делать выводы из фактических доказательств, и потраченное на формулировку вопроса время поможет нам получить точный ответ, не зависящий от собственных и чужих предположений.

Это первый этап процесса анализа данных. Аналитики данных должны проявлять здесь некоторую креативность. Мы не меняем информацию в соответствии с нашими идеями, мы формулируем идеи, чтобы добиться полезного для нас понимания. В главе 4 «Сформулируйте вопрос» мы исследуем различные методы, а их применение обеспечит соответствие вопросов, которые мы в конечном итоге зададим нашим данным, целям проекта и удержит нас от пропусков и «расползания границ проекта» — неконтролируемого выхода проекта за первоначально установленные рамки условий.

Правильные ингредиенты

Мы уже давно вступили в эру компьютеров, и большинство учреждений государственного и частного секторов накопили огромное количество своих собственных данных. Однако данные собирались задолго до того, как мы узнали, что с ними можно делать, и зачастую это делали сотрудники, которые не знали, как исследовать, стандартизировать и анализировать информацию, чтобы она действительно была полезной. Такой пробел в знаниях способен вызвать в лучшем случае организованный хаос, когда массивы данных могут содержать искаженные и грязные данные, о которых мы узнаем больше в главе 5 «Подготовка данных».

Если вам надо очистить данные и сделать их удобочитаемыми, нельзя торопиться. Чтобы понять, насколько важно подготовить данные, прежде чем делать с ними что-либо, обратимся к процессу оптимального распознавания символов (OCR) при сканировании. Программное обеспечение OCR отсканирует страницу письменного или печатного текста и переведет этот текст в цифровой формат. Но OCR-сканы не всегда на 100% корректны: их точность зависит как от возможностей программного обеспечения, так и от качества распечатываемой страницы. Рукописные документы XVII в. создадут больше трудностей и спровоцируют больше ошибок, которые затем должны быть вручную исправлены в более поздних данных. Те, кто не знает, как правильно записывать данные, или кто использует установленные в учреждении устаревшие или неоптимальные стандарты, будут генерировать массивы данных, которые также должны быть «очищены».

Игра в действии

Для анализа современных данных не требуется такой же уровень осторожности, как на предыдущих двух этапах. Если вы нашли время на формулирование правильного вопроса и подготовку своих данных для того, чтобы уяснить, что от них требуется, вы можете позволить себе поэкспериментировать с анализом. Прелесть работы с массивами данных заключается в том, что вы можете дублировать их, поэтому работа с одним типом алгоритма на массиве данных не исключает

возможности применения к нему и другого алгоритма. Этим хороша цифровая информация — ее можно использовать, отбирать, реструктурировать и извлекать, но вы все равно можете вернуться к более ранней версии, как только закончите работу, и начать снова.

Итак, вы потратили время на создание лесов для вашего проекта и обеспечение того, чтобы они не рухнули под тяжестью вопросов, которые вы задаете, так что теперь пришло время исследования. В главах 6 и 7 приведены решения для типов анализов, которые вы можете выполнять, а также краткий перечень их преимуществ и ограничений, чтобы повысить вашу уверенность в выборе алгоритма, оптимального для целей конкретного проекта.

Начало работы

Хотя эта часть в основном теоретическая, она имеет практическое значение, и поэтому я настоятельно рекомендую рассмотреть возможность применения каждого из пяти этапов, описанных выше, к вашему собственному проекту параллельно с чтением книги. Тогда вы освоите некоторые из необходимых инструментов, прежде чем начать изучение этой части.

Массив данных

Если у вас еще нет собственного массива данных, с которым вы можете работать, не волнуйтесь. Существует множество общедоступных массивов данных — вы можете бесплатно использовать их в собственных экспериментах. Большим преимуществом является то, что вы сразу же погрузитесь в использование реальных массивов данных, а не тех, что были специально созданы для обучения. По моему опыту, реальные массивы данных позволят вам испытать чувство победы в результате извлечения идей из реальной информации, и добавят вес утверждению, что наука о данных имеет важное значение для будущего развития огромного количества дисциплин.

Действительно интересных и разнообразных массивов данных, доступных в интернете для загрузки и использования, очень много, однако выбор за вами. Вот только несколько для начала:

- **World Bank Data.** Данные Всемирного банка — ценный ресурс глобальных данных о развитии.
- **European Union Open Data Portal.** Портал открытых данных Европейского союза — правительственные данные государств — членов ЕС.
- **Million Song Dataset.** Сборник метаданных и аудиозаписей популярной музыки.
- **The CIA World Factbook.** Всемирный справочник ЦРУ — массивы данных из 267 стран по темам от истории до инфраструктуры.
- **National Climatic Data Center.** Национальный центр климатических данных — сведения об окружающей среде США.

Программное обеспечение

Новичку в науке о данных необходимо понять, что данные не имеют своего собственного «языка» и что они могут «говорить» с нами только через машину или элемент программного обеспечения. Под «языком» данных я здесь подразумеваю способ, которым машина передает данные специалисту. О скорости автомобиля, цикле цветения растений, температуре наружного воздуха и количестве жителей в городе данные просто *есть*. Это ряд компонентов, но отношения, которые мы устанавливаем между ними, — дело человека или компьютера. Если продолжить аналогию с языком, я бы сказал, что данные можно сравнить с отдельными буквами, которые ждут, когда кто-то расставит их в соответствующем порядке, чтобы сформировать слова и предложения. Таким образом, от нас зависит (через инструменты, которые мы применяем), чтобы наши данные работали.

Наличие доступа к программному обеспечению не является обязательным требованием для тех, кто изучает данные с помощью этой книги, поскольку она сосредоточена на практическом применении,

а не на кодировании. Но если вы хотите попробовать использовать некоторые из приведенных здесь примеров, я рекомендовал бы либо R, либо Python — оба этих языка представляют собой программные средства анализа данных и доступны для бесплатной загрузки в Windows, Linux/Unix и Mac OS X. В настоящее время это два самых распространенных в отрасли инструмента с открытым исходным кодом.

Сформулируйте вопрос

04

Я часто слышу, как другие аналитики данных сетуют на то, что данных *слишком много* и что сама идея разобраться с таким количеством информации для ответа на бизнес-вопрос ошеломляет. С учетом почти постоянного потока «выхлопных данных» как мы можем надеяться управлять собранной информацией таким образом, чтобы это способствовало ее рассмотрению? Мы не можем просто прогнать все имеющиеся у нас сведения через некий алгоритм и скрестить пальцы в надежде получить нужные нам результаты.

Прежде чем мы сможем подготовить и проанализировать данные, мы должны знать, сведения какого рода нам нужны. А для этого необходима небольшая тонкая настройка вопросов нашего проекта.

Руководители часто ставят проблему перед аналитиком данных и ожидают, что тот сразу же погрузится прямо в базу данных. Но сначала поставленный вопрос нужно понять, разобрать, проанализировать. Мы должны знать, о чем нас спрашивают; если мы не ответим должным образом на этот вопрос, результаты проекта будут бесполезны. Рассмотрим процесс написания школьной работы: действие наобум, попытка ответить на вопрос, как только вам его задали, приведет (если только вы не суперсчастливчик) к тому, что вы просто дадите кучу громоздкой, неструктурированной информации. Только если вы найдете время, чтобы сделать шаг назад и подумать о картине в целом — рассмотреть ее многочисленные компоненты и контекст, — можно будет говорить об убедительности и логичности ваших аргументов.

Полное понимание вопроса также помогает придерживаться курса и снижает риск отклонения от поставленной цели. Допустим, наш учитель истории хотел, чтобы мы написали об американской войне за независимость. Анекдоты из биографии Джорджа Вашингтона могут

быть интересны, но они не отвечают заданной теме. Эссе с такой несоответствующей информацией созданы учениками, которые погружаются в предмет, не поняв сути вопроса, и вместо этого используют все, что они могут собрать, не заботясь о том, чтобы отказаться от ненужных данных.

Именно поэтому *в первую очередь* нам необходимо определиться с вопросом.

В главе 4 я покажу наиболее подходящий/плодотворный способ действий на этой стадии процесса анализа данных. Поскольку определение вопроса может показаться чрезвычайно сложной задачей, я предлагаю вам проверенный на практике подход, который проведет вас через этот этап и обеспечит рассмотрение всех аспектов вопроса, а также защитит вас от боссов, которые стремятся навязать дополнительную работу после начала проекта.

Смотри, мама, никаких данных!

Несмотря на всю важность, выявление проблемы, как правило, является наиболее часто игнорируемой частью проектов, использующих данные. Я и сам грешил этим, поскольку долгое время начинал проекты с подготовки данных. Но это было не потому, что хотел проскочить вперед; я просто думал, что постановки проблемы достаточно. В конце концов, аналитики данных часто привлекаются к работе над *проблемами*, и на моей первой работе в Deloitte все проекты начинались с технического задания, в котором уточнялось то, что от меня требуется, и указывалось, где нужна помощь. Неудивительно, что фирма стандартизировала и упорядочила процесс, но это только сформировало во мне плохую привычку бежать впереди паровоза, прежде чем представить себе целостную картину.

Еще одна причина пренебрежения определением вопроса заключается в том, что на этом этапе не используется много данных (если они вообще используются), в результате чего многие аналитики данных относятся немного снисходительно к выполнению этого этапа. Но важно отметить, что те, кто предложил вопрос, вероятно, не являются специалистами по данным и не знают о подготовке, необходимой

для очистки и анализа данных. Немногие компании на сегодняшний день информируют своих сотрудников о важности хранения информации и обеспечения доступа к ней, и этот пробел в знаниях означает, что у многих аналитиков данных все еще спрашивают: «У нас есть много данных, может быть, вы сформулируете какие-то идеи на их основе?» Подобные вопросы задаются очень часто, хотя они туманны, расплывчаты и ничему не способствуют в процессе решения проблемы компании*.

Таким образом, даже если вопрос был сформулирован заранее и ваш босс спрашивает, почему вы не занимаетесь обработкой данных, не думайте о себе как о сумасброде. Изложите свои доводы. Просто предложить вопрос недостаточно — он должен быть переформулирован в терминах, которые будут соответствовать имеющимся данным, иначе реализация проекта застопорится.

Как решить такую проблему, как...

Задания, которые поступают от руководства организации или от инвесторов, часто постулируются как открытые пути к вопросу, а не реальный вопрос как таковой: «Мы недопоставляем единицы продукции», или «Наши клиенты покидают нас быстрее, чем ожидалось», или «В нашем продукте есть дефект». Ни одно из этих утверждений не является вопросом. Я призываю читателей применять следующий поэтапный подход к выявлению и решению проблемы на основе данных. Это сделает первый этап более эффективным и снизит риск того, что вы сосредоточитесь на неправильной проблеме.

1. Поймите проблему

Любой, кто планирует участвовать в проектах, связанных с данными, должен прежде всего знать о ловушке, в которую он может невольно попасть еще до того, как работа начнется: если последовать решению

* Многие руководители считают, что для выявления проблемы следует использовать данные, однако такой подход редко срабатывает. Мы не можем заставить данные говорить — мы должны сначала определить, что хотим услышать.

коллеги относительно того, какими вопросами следует заняться, можно фактически выбрать ошибочную проблему. Коллега может иметь благие намерения и пытаться быть более полезным, представляя разработанные им вопросы, но они необязательно будут пригодны для запроса, адресованного данным. Заманчиво, получив несколько, по-видимому, хорошо сформулированных запросов, не беспокоиться о том, чтобы идентифицировать вопрос самостоятельно. Но это может привести к катастрофе в дальнейшем процессе; именно от вас зависит определение всех параметров бизнес-проблемы, потому что вы *обучены* тому, как это сделать. Слепое заимствование набора вопросов у специалиста, не связанного с данными, и применение их к вашему проекту может увенчаться решением не той проблемы или просто не привести к каким-либо результатам, потому что у вас нет для них необходимых данных.

Прежде чем начать непосредственно трудиться над проектом, в первую очередь нужно *поговорить* с человеком, который поставил перед нами проблему. Понимание не только того, *что* это за проблема, но и почему она должна быть решена *сейчас*, кто основные заинтересованные стороны и что ее решение будет *означать* для учреждения,

Хватит нюхать розы

Несмотря на то что хороший аналитик данных должен начать с постановки проблемы и иметь четкое представление о ее различных аспектах, не беспокойтесь, если на этом этапе проявятся несколько неожиданных отклонений. На первой стадии процесса мы активно изучаем архив информации компании, большая часть которого, возможно, никогда ранее не обрабатывалась, — а также проводим весомую черновую работу, чтобы понять проблему компании, поэтому нас не должны пугать сюрпризы на пути. Не списывайте эти отклонения как второстепенные, отмечайте их по мере своего продвижения вперед. Именно этим маршрутом и является наука о данных, а дополнительная аналитика — ее важная часть. Иногда отклонения могут не иметь отношения к проблеме, но обладать дополнительной ценностью для бизнеса. Порой я обнаруживал, что внезапные идеи могут принести еще большую ценность, чем само решение проблемы, которой меня попросили заняться.

поможет начать «отладку» нашего исследования. Без этого шага результат может оказаться опасным для аналитика данных, так как в дальнейшем по ходу проекта мы, несомненно, интерпретируем поставленный вопрос иначе, чем заинтересованные стороны. Как только мы разобрались с центральной проблемой, можно перейти ко второму шагу.

2. Изучить отрасль

Если у вас уже имеются знания о сфере, в которой вас просят работать, это отличное начало. Вы можете применить свой опыт. Вы, возможно, уже знаете, например, конкретные проблемы, с которыми обычно сталкиваются компании, работающие в этом секторе, или можете быть в курсе того, какие отделы, как правило, занимались этими вопросами успешно или тщетно либо какие конкурирующие компании обнаружили и решили именно те проблемы, которые были поставлены перед вами.

Если у вас нет знаний об отрасли, не все потеряно. Потратьте некоторое время, исследуя ее более подробно. Каковы подводные камни в отрасли? Столкнулись ли конкуренты вашей компании с аналогичными проблемами или же есть существенные различия? Как они решали подобные проблемы? Миссия и цели компании, для которой вы работаете, существенны для отрасли как таковой? Чем эта компания отличается от других по объему производства, организационной структуре и рабочему процессу?

Google может быть вашим лучшим другом в поиске ответа на многие из этих вопросов, но также помните, что как аналитик данных вы не работаете в вакууме. Детальное знание среды, где вы действуете, а также ее индивидуальных особенностей и ограничений поможет вам разработать подход, значимый для тех, кто так или иначе связан с проектом. Не будьте отшельником. Если у вас пробелы в знаниях, используйте свой лучший ресурс — коллег. И даже если у вас есть вся необходимая информация, все равно пойдите и поговорите о том, что вы узнали, с соответствующими сотрудниками. Люди, которые вызвали ваш проект к жизни, всегда будут хорошей отправной точкой, чтобы убедиться, что вы говорите с теми, с кем нужно. Они не только помогут вам узнать недостающие сведения, но и направят к должностным лицам, ответственным за те участки в организации, где возникла проблема.

Будьте общительным человеком

Проекты, основанные на работе с данными, часто затрагивают не одно направление деятельности компании. Вы обнаружите, что для того, чтобы действительно решить вопрос, вы должны контактировать с несколькими подразделениями, а также с людьми, которые часто не имеют оснований общаться друг с другом. Будьте готовы к коммуникации, одной из многих радостей аналитика данных. В конце концов, доступ к информации обо всей организации ставит вас в уникальное положение, позволяющее собрать как можно больше сведений о том, как она работает, и о том, как бизнес-вопрос, для решения которого вы были призваны, повлияет на деятельность всей компании. Это весьма захватывающе.

Однажды, когда я анализировал для некой организации опыт ее работы с клиентами, я был технически связан с командой отдела маркетинга по работе с клиентами и аналитике. Но половину времени фактически потратил на работу в IT-отделе: естественно, мне нужен был доступ к их базам данных, серверам, инструментам и программному обеспечению. До того, как я пришел, в компании не было аналитика данных и эти два отдела практически не пересекались, редко находя причины общаться.

Начиная проект, вы увидите, что можете послужить мостиком между отделами, помогая синхронизировать их усилия в части данных. Это поспособствует изучению того, как компания функционирует и как ваш проект способен усовершенствовать их услуги и методы работы. К концу моего пребывания в той организации я работал с отделами маркетинга, IT и отделом оперативного управления, а мои отчеты часто собирали их всех за одним столом для обсуждения планов на ближайшие недели.

Каждый раз, когда вы контактируете с новым отделом, убедитесь, что вы выстраиваете связи так, чтобы команда знала о работе, которую вы делаете. Тем самым вы будете сохранять канал связи открытым на случай, если вам в дальнейшем понадобится информация от этих сотрудников.

3. *Думайте как консультант*

Большинство согласится с тем, что наука о данных требует подхода «снизу вверх»: мы используем данные компании для проведения анализа и постепенно выстраиваем на них наши результаты, чтобы лучше понять внутренние проблемы компании. Этот потенциал данных и является тем, что делает первый шаг настолько захватывающим. Но для того, чтобы выявить вопрос, нужно обращать больше внимания на методы бизнес-консалтинга.

В консалтинге мы выделяем возможные стратегические подходы для бизнеса. В качестве консультантов, как правило, выступают люди, проработавшие в бизнесе или отрасли несколько лет и накопившие много знаний о соответствующем секторе. Они часто занимаются улучшением крупномасштабных стратегических и организационных аспектов компании, что требует нисходящего подхода, — и такая методология анализа общей картины обязывает делать определенные предположения о поставленной проблеме.

Для нас может показаться контрпродуктивным использовать консалтинговые методы; как аналитикам данных нам советуют стараться воздерживаться от предположений и использовать как можно больше достоверных данных. Но пример консультантов может быть чрезвычайно полезен, особенно на первом этапе процесса. Итак, забудьте на мгновение о своих технических знаниях и посмотрите на организацию, участников проекта и стратегию компании, прежде чем начать размышлять о данных. В конце концов, определение вопроса касается фильтрации имеющихся вариантов, и этот третий шаг поможет уточнить вопросы, чтобы они стали соответствовать потребностям нашей компании.

Составьте список ключевых участников проекта и особо отметьте, кто будет принимать окончательное решение. Проведите с ними некоторое время и не переходите к четвертому шагу, пока не сможете ответить на следующие вопросы:

- Что каждый из участников проекта думает о проблеме?
- Каковы грани этой проблемы?
- Какие отделы должны быть в фокусе моего внимания?
- В чем могут быть первопричины проблемы?

- Считают ли участники, что я должен поговорить с кем-либо еще? Я с ними разговаривал?
- Где находятся данные и кто за них отвечает?
- Что будет означать успех этого проекта?

Сверху вниз или снизу вверх?

Deloitte, PricewaterhouseCoopers, KPMG, Ernst & Young — известные компании «Большой четверки», которые применяют консалтинговый подход «сверху вниз» ко всем своим проектам. Но что такое нисходящий метод, чем он отличается от восходящего и действительно ли один лучше другого?

Давайте рассмотрим эти вопросы один за другим. Нисходящий подход предполагает, что сначала изучается проект в целом, а затем — его детали. Консультанты, привлеченные на новый проект, поначалу действуют на верхних ступенях иерархии компании: читают соответствующие отчеты и затем последовательно переходят от контактов с гендиректором к исполнительным директорам, менеджерам и т.д. — пока не достигнут людей «внизу» (если вообще туда доберутся). Какими бы тщательными ни были их исследования, консультанты, применяющие нисходящий подход, в конечном итоге будут представлять свой проект на основе первоначальных выводов, подобных тем, что приведены в финансовых отчетах компании на конец года. Нисходящий подход остается типичным методом, используемым консалтинговыми фирмами.

Восходящий подход — прямо противоположный путь. Этот метод на *первое место* ставит *цифры*. Сначала рассматриваются данные компании, и только потом начинается движение вверх по иерархической лестнице, от менеджеров проектов до руководителей подразделений. Количество собранной информации нарастает — и только потом можно обращаться к руководителям высшего звена. Этот подход опирается на факты и цифры, сообщающие исследователю о повседневной работе компании. Читатели не удивятся, узнав, что восходящий подход почти всегда используется аналитиками данных.

Учитывая, что это книга о науке о данных, мы могли бы автоматически взять на вооружение восходящий подход, думая, что метод, который фокусируется на базах данных, безусловно, предпочтительнее, чем тот, который начинается с отфильтрованных данных. Это, однако, было бы опережением событий. Да, поскольку подход «снизу вверх» основан на фактах, он позволяет прийти к выводу гораздо быстрее, чем если бы мы двигались «сверху вниз». Тем не менее любой исследователь проекта, использующий восходящий подход, скажет вам, что с помощью этого метода практически невозможно провести изменения. Компании не просто работают с помощью данных — они работают с людьми и отношениями. Мы не можем вгрызаться в цифры, а затем ворваться в кабинет директора, заявляя, что знаем, как решить проблему, размахивая бумажкой с данными, чтобы доказать это. Цифры, как вы узнаете из этой книги, всего лишь одна часть головоломки. Мы также должны понимать культуру компании, ее миссию и сотрудников.

Поэтому, если мы хотим прийти к достоверным выводам, нужно еще немного детективной работы. Для этого идеально подходит нисходящий подход. Ричард Хопкинс, мой наставник и директор PricewaterhouseCoopers, говорит, что, если вы только следуете за цифрами,

«компания не будет вас слушать, потому что вы не сотрудничали и не разговаривали с ними. Предположим, вы выяснили на основе чисел, что продажа определенного продукта приводит к потере денег, и делаете вывод, что они должны прекратить его производство. Но они могли продавать такие продукты по определенной причине — например, для того, чтобы продать нечто другое; или, возможно, выпуск этого продукта уже сворачивается. Да, подход “снизу вверх” быстро приведет вас к результату, на основе которого вы сможете принять решение. Но без понимания общей картины это решение рискует оказаться нелучшим».

(SuperDataScience, 2016)

Нисходящий подход обеспечивает вовлеченность всех участников. Мы делаем выводы благодаря вкладу коллег. С таким уровнем сотрудничества мы сможем добиться изменений в компании гораздо быстрее.

Итак, вместо того чтобы использовать метод «снизу вверх», следует ли нам подражать сотрудникам KPMG и использовать подход «сверху вниз»? Не обязательно. На то, чтобы понять общую картину, требуется много времени. Мы должны работать с участниками проекта, организовывать семинары — и когда нам дают так много информации на столь раннем этапе, мы вынуждены принимать важные решения о том, какая информация полезна, а какая нет. Это большая ответственность. Более того, результаты, полученные по итогам использования нисходящего подхода, могут в конце концов не соответствовать данным.

Что делать? Ричард рассказал о явных преимуществах применения обоих этих методов к одному и тому же проекту:

«Довольно часто возникает разрыв между числами, полученными из баз данных, и числами, взятыми из заключительных отчетов. Я обнаружил, что достигаю оптимальных результатов, когда провожу оба анализа, а затем объединяю их. Это позволяет понять, что происходит с итоговыми отчетами и какие факторы приводят к этим финансовым результатам. Таким образом, использование подходов “сверху вниз” и “снизу вверх” связывает данные с процессом и дает нам полную картину».

(SuperDataScience, 2016)

Лучше всего, если возможно совместить два подхода. Благодаря этому вы сможете не только быстро получить ответ с необходимой информацией, но и заручиться поддержкой участников проекта, столь необходимой для того, чтобы реализовать свои плодотворные идеи.

4. Осознайте ограничения

Что делать, если после всей проделанной здесь работы мы обнаружим, что данных, которые мы расценили как необходимые для нашего исследования, нет?

Наиболее эффективный подход заключается в том, чтобы изучить высший уровень массива данных и понять, какие данные у нас на самом деле есть и необходим ли дополнительный сбор сведений, прежде чем проект сможет двигаться вперед. Опять же, это означает общение с правильными людьми — теми, кто отвечает за существующие данные компании. Благодаря такому общению мы сможем получить более полное представление о том, где в данных могут обнаружиться проблемы и где нам может понадобиться дополнительная информация, чтобы гарантировать статистическую значимость выбранных нами данных. Этот шаг немного напоминает дилемму курицы и яйца, ведь мы должны знать, какие вопросы задать данным, прежде чем мы увидим эти базы данных. Но мы также должны убедиться в том, что уже на ранней стадии имеем нужные данные, — иначе потеряем много времени, прежде чем приступим к следующему этапу процесса анализа данных.

Практика — лучший путь к освоению этого. Вспомните свои мысли о том, какие типы данных будут полезны для ответа на ваши вопросы. Напишите их рядом с вопросами и делайте отметки, чтобы понять, что вам нужно в каждой контрольной точке. На первом проекте это может напомнить одновременное жонглирование многими тарелками, но с опытом все станет намного проще.

Кейс: восполнение пробелов

Рубен Коугел — руководитель отдела данных калифорнийской технологической компании VSCO, фокусирующейся на сфере искусства. На базе онлайн-подписки компания дает художникам возможность создавать цифровые инструменты взаимодействия с пользователями. На момент своего назначения Коугел был в VSCO единственным аналитиком данных, и стандартизированная практика создания отчетов, основанных на данных, в компании отсутствовала. Но Коугел увидел в имеющихся данных возможность ответить на ключевые вопросы, важные для компании: кем являются люди, которые покупают подписки VSCO, и ведут ли они себя «по-другому» после покупки?

Рубен знал, что проблема требовала обращения к бесплатным учетным записям пользователей, перешедших впоследствии на платные услуги подписки. Но это была только верхушка айсберга — Рубену нужно было начинать «бурение» и копать все глубже:

«Мне требовалось больше информации для того, чтобы начать работу. Я хотел иметь представление еще и о том, что именно указывает на изменения в поведении и почему важно “знать” наших клиентов. В то время мне не было известно, каким образом VSCO выявляет свою целевую аудиторию, так что это был пробел в знаниях, который мне следовало восполнить, прежде чем я мог бы найти ответы».

(SuperDataScience, 2017)

Рассмотрев проблему с такой точки зрения, Рубен смог определить фокус анализа. В ходе его бесед с коллегами обнаружилось, что проблема в действительности связана с маркетингом. Таким образом, вопрос был поставлен по-новому — с учетом целевых потребностей маркетинга: «Миллионы пользователей VSCO являются потенциальными покупателями, но не все они одинаково склонны покупать подписку VSCO. Тогда выясним: 1) как сегментируются наши пользователи с точки зрения их предпочтений, поведения и демографии и 2) какие из этих клиентских сегментов представляют наиболее вероятных покупателей?»

После того как это прояснилось, остальное стало очевидно для Рубена. От сотрудников отдела маркетинга он знал, какое поведение предсказывало, что клиент, вероятнее всего, приобретет подписку. Рубен эффективно преобразовал изначально довольно неструктурированный вопрос в нечто, что не только обеспечило инвестора в точности той информацией, которая ему требовалась, но и указывало, где компания может собирать больше данных для того, чтобы улучшить свою маркетинговую практику.

Если вы работаете в компании, которая с течением времени накопила много данных, первоначальное выявление проблемы становится еще более важным делом. Я часто обнаруживал, что, хотя коллеги могут положительно воспринимать идею использования данных, они не вполне уверены, что данные могут им помочь. Это может относиться даже к людям, управляющим базами данных. И именно здесь многие компании делают неверный шаг. По иронии судьбы из-за того, что так велико количество данных, имеющихся в распоряжении у компаний, последние перестают осознавать смысл и значение данных и, следовательно, преуменьшают их ценность. Мы должны помочь им понять важность данных, и это путешествие начинается здесь.

Заручитесь согласием других участников

Ваш анализ данных *встряхнет* обстановку. Как возмутитель спокойствия, вы можете столкнуться с сопротивлением. Если вам намеренно мешает коллега, примите меры. Если вы упускаете информацию в результате намеренного противодействия сотрудника, не стесняйтесь заострить этот вопрос. Суть любого проекта в области науки о данных в том, чтобы повысить ценность компании, и, если заинтересованные стороны знают об этом, они также должны знать, что данные приоритетны. Вы не склоните всех к этой идее, но я уверен, что нельзя сдаваться: будьте готовы преодолевать сопротивление, чтобы выполнить свою работу.

5. Проведите майнинг данных (по желанию)

Глубинный анализ (майнинг) данных — возможно, самая приятная для меня часть процесса в любом проекте. То, что ученым не дают заниматься майнингом данных, немного похоже на запрет кураторам музеев изучать материалы, за которые они несут ответственность. Именно на этом этапе вы можете позволить себе быть исследователем. Для меня глубинный анализ данных — процесс, в котором вы выполняете тестирование с целью тщательного изучения данных на высшем уровне и находите области, которые могут предложить интересные идеи для дальнейшего исследования. На этом экспериментальном этапе мне нравится помещать данные в Tableau*, которое умеет их читать и поможет вам создать предварительные наглядные визуализации, такие как легко читаемые таблицы, диаграммы и графики. Это обеспечивает прекрасный задел, который вы можете использовать в качестве фокусирующей линзы, чтобы сформулировать нужные вопросы.

В конечном счете, если майнинг данных выполняется на начальном этапе проекта, он наиболее эффективно помогает лучше понять

* Программное обеспечение для визуализации, которое я буду обсуждать более подробно в главе 8 «Визуализация данных».

проблему и управлять процессом анализа. Это тест-драйв ваших данных: вы испытываете их в необработанном виде, чтобы увидеть, могут ли какие-либо тенденции проявиться даже на раннем этапе. Майнинг данных может сэкономить много усилий в дальнейшем. В то же время не унывайте, если он ни к чему не приведет. Данные могут предлагать или не предлагать нам дальнейшие действия или решения в зависимости от нескольких факторов, таких как компания, качество данных и уровень сложности проблемы. Итак, сделайте этот шаг, но не забывайте «делить на десять». И если вы найдете что-то интересное, запишите и убедитесь, что вы помните о своих находках, когда перейдете к шестому шагу...

6. Уточните проблему

Теперь, когда мы поняли масштаб проблемы и определили количество данных, имеющихся в нашем распоряжении, можно начать копать немного глубже. Здесь мы начинаем сопоставлять масштаб проекта с данными, чтобы отделить переменные и данные, которые будут полезны, от тех, которые не пригодятся, и чтобы надлежащим образом переформулировать вопрос.

Хотя все данные потенциально могут оказаться полезными, мы не можем использовать всю имеющуюся информацию по каждой проблеме, и это только к лучшему: если бы все данные были полезны всегда, объем получаемой на выходе информации был бы просто слишком громоздким для управления. По этой причине мы можем быть разборчивыми по отношению к предоставленным данным. Это означает, что мы должны учесть параметры и контекст проблемы, которую хотим решить, прежде чем двигаться вперед. В конечном счете уточнение проблемы экономит время, устраняя данные, которые не имеют отношения к нашему вопросу.

7. Соберите дополнительные данные

На этом этапе вы уже определили, какие данные вам нужны, и составили разумный перечень вопросов для решения проблемы. Сейчас самое время оценить эффективность ваших подвопросов. В конце концов,

Разделение данных

Давайте возьмем, к примеру, компанию с низкой прибылью. Ее главная проблема — недостаточная рентабельность. Возможно, нам сразу захочется спросить: каковы причины этого? Возможно, общая *выручка* компании слишком мала. Или *расходы* компании слишком высоки. Это уже две возможности. Используя эти две категории (доходы и расходы), мы сможем дополнительно проанализировать проблему. Какие расходы несет компания? Какие различные продукты и услуги она предлагает? Задавая эти дополнительные вопросы, мы можем обратиться к заинтересованным сторонам с конкретными запросами: вам понадобится ознакомиться с финансами компании, увидеть, как сегментировано предложение товаров/услуг, определить, какие цели по прибыли компания ставила себе на протяжении времени, и т.д.

Если участникам проекта известны эти сведения, они могут дать прямой ответ на некоторые из ваших вопросов. Например, вполне вероятно, что они будут знать, существенно ли изменились расходы за определенный период времени. Просто задав этот вопрос, мы можем исключить проблему расходов, что позволяет нам потенциально* свести проблему к одним доходам.

просто не стоит отвечать на те вопросы, которые, как вы поняли, компании не интересуют или по которым ничего не будет предпринято. Спросите себя сейчас: каковы ожидаемые результаты от этих подвопросов? Помогают ли они решить проблему или чего-то еще не хватает?

Именно здесь вы благодарите себя за то, что прошли предыдущие шесть шагов, достигнув этой точки; выделение ключевых областей, из которых вам нужны дополнительные данные, оптимизирует и, следовательно, ускорит процесс сбора данных. Составьте план, а затем отложите его в сторону; воздержитесь от сбора каких-либо данных вообще, пока вы не выполните восьмой шаг.

* Я говорю «потенциально», потому что важно не сбрасывать со счетов другие варианты слишком рано в этом процессе. Следите за решениями, которые вы принимаете, и запишите, как и почему вы посчитали массив данных менее важным. Это позволит вам быстро вернуться к началу, если позже нужно будет переформулировать вопрос.

Количественные и качественные методы

Если мы определили, что нам нужно больше данных (и это очень вероятно; я не участвовал ни в одном проекте, который не потребовал бы дополнительных данных), именно на этом этапе мы должны рассмотреть, данные *какого рода* обусловят наилучшие результаты. Наиболее важные вопросы для сбора данных: «Где *находятся* источники?» и «Мы хотим использовать количественные или качественные методы исследования?».

Что такое количественные и качественные методы? Проще говоря, с помощью количественных методов собирают числовую информацию, в то время как посредством качественных — нечисловую. Но есть несколько дополнительных отличий, которые следует принять к сведению, прежде чем определить, какой метод использовать при сборе дополнительных данных.

Количественные методы

Количественные методы следует применять, когда нужна статистическая информация. Результаты их применения, как мы увидим в следующей главе, *гораздо* проще собрать в массив данных, чем использовать качественные методы, но наше окончательное решение о том, какие данные собирать, не должно основываться на принципе «простота ради простоты». Каждый метод отвечает на свои вопросы. Мы не можем, например, предпочесть качественный подход к сбору сведений о возрасте, потому что возраст — это факт, а не мнение (что бы вы себе ни говорили). Мы будем использовать количественные методы, если нужно подсчитать элементы, или измерить изменения в заработной плате, или узнать больше о демографии потребителей. Не забудьте отметить, что количественные данные не только числовые сведения; скорее, это данные, которые могут быть подсчитаны. Вопросы о любимых брендах респондентов или политической принадлежности все еще предполагают использование количественных методов, потому что ответы технически могут быть подсчитаны по категориям.

Качественные методы

Качественные методы связаны с открытыми вопросами, имеющими бесконечное количество ответов. По своей природе они носят исследовательский характер и помогают выявить — но не количественно — тенденции во мнениях, мыслях и чувствах. Мы можем применять этот подход, когда нужно больше контекста, чтобы понять проблему, или когда проблема слишком сложна, чтобы решить ее количественным методом. Таким образом, качественные методы лучше подходят для сбора данных об эмоциональных инвестициях потребителя в продукт или тогда, когда нужно дать более развернутый ответ на вопрос, как респонденты могут относиться к какой-либо политической партии.

8. Проинформируйте заинтересованные стороны*

После того как мы приняли во внимание все предыдущие семь шагов, крайне необходимо, чтобы у нас, нашей команды и *всех участников проекта* было общее понимание ситуации. Четкая и ясная постановка проблемы, которую вы будете решать, обеспечит точный выбор подхода, и это уменьшит шансы других изменить ориентиры в процессе реализации проекта.

Сторона, обратившаяся с просьбой о выполнении проекта, *должна согласиться* с вашим планом решения проблемы, который в идеале должен включать в себя не только то, что касается содержания проекта, но и его временные рамки. Я настоятельно рекомендую разделить проект на этапы, что позволит всем вовлеченным лицам оставаться в курсе вашего продвижения вперед и защитит вас от любой негативной реакции в конце проекта и упреков в том, что вы скрывали свои намерения.

Также необходимо объяснить заинтересованным сторонам, что это *не* обычный бизнес-проект, что проекты в области науки о данных

* Возможно, данное действие не поможет определить параметры, необходимые для постановки вопроса, но тем не менее крайне важно, чтобы вы выполнили этот этап.

не всегда соответствуют моделям PRINCE2, которые столь знакомы и любимы бизнесом. Это поможет защитить вас от предвзятого вмешательства и даст возможность точно объяснить участникам, какие шаги вы *собираетесь* предпринять для выполнения задачи.

Единственное, на чем я настаиваю в начале любого проекта, связанного с наукой о данных, — это *письменное подтверждение согласия заинтересованных лиц*. Вы можете быть лучшими друзьями в личной жизни, но по моему опыту участники, в каком бы качестве они ни выступали, в ходе реализации проекта склонны менять свое представление о том, чего они хотят. Такое поведение понятно в случае, когда сам проект имеет расплывчатый характер, но оно способствует разрастанию масштаба, которое может либо вывести вас за рамки исходных параметров, либо полностью убить проект. Поэтому, прежде чем перейти к подготовке данных, получите письменное подтверждение согласия.

Правда не всегда приятна

Как справиться с ситуацией и не оттолкнуть других вовлеченных в проект людей, если ваши результаты не дают им того, на что они надеялись? Как избежать этого минного поля? К сожалению, вы не можете сделать ставку на то, что говорят вам данные, если правда является негативным результатом. Правда не всегда приятна. В науке о данных вы смотрите на неопровержимые факты — игнорирование или приукрашивание результатов ради спасения чувств участников проекта ни к чему не приведет.

Если я могу дать здесь какой-либо совет, вот он: подготовьте людей, с которыми имеете дело, к возможности того, что результаты могут быть не такими, как им хочется. Поясните заранее, что вы не знаете, каковы будут результаты. Заказчикам могут не понравиться итоги. Сразу же дайте понять, что вы добываете факты, а не стремитесь польстить, и надеетесь, что они не будут «убивать гонца», если отчет окажется не особенно благоприятным.

Соблюдение графика

Этот этап процесса анализа данных не должен продолжаться чересчур долго в цикле проекта. Иногда новички могут потратить на него слишком много времени, потому что хотят убедиться, что они разработали надежную методологию. Помните: вы никогда не сможете прояснить проблему до такой степени, чтобы точно знать, чего хотите. Если вы проделали хорошую работу на этом этапе, то, скорее всего, сэкономите время, но вы также должны научиться позволять процессу идти своим чередом — это умение приходит с опытом.

Если вы будете последовательно выполнять шаги, описанные выше, это застрахует вас от дальнейших трудностей и поможет обрести уверенность для перехода к своевременной подготовке данных. В конечном счете, если задача, которую перед вами поставили, не является дьявольски сложной и не требует многочисленных согласований, выявление и уточнение проблемы должно занять максимум неделю. Но старайтесь по возможности не ставить других в известность о сроках, которые вы стремитесь соблюсти, — это только добавит давления на вас. Если для вашего комфорта и прогресса требуется еще несколько дней, тем лучше.

Моя рекомендация? Дайте себе достаточно времени, чтобы уложиться в срок. Гораздо лучше пообещать меньше и перевыполнить обещание, чем сделать обратное. Полезно сначала определить, сколько дней, по вашему мнению, займет проект в целом, а затем добавить 20% к этому количеству. Чаще всего на анализ данных времени не хватает. И если вы столкнетесь с какими-либо препятствиями и подумаете, что не успеете завершить работу к дате, о которой вы изначально договорились, не забудьте предупредить заказчика — он должен узнать об этом как можно раньше. Информирование людей укрепит доверие между вами и другими участниками проекта и сделает их вашими единомышленниками.

Искусство говорить «нет»

По мере того как будете укреплять свои позиции в компании, приумножать ее успехи и повышать эффективность сотрудников, люди будут

чаще приходиться к вам за помощью в их проектах. К сожалению, у вас не будет времени помочь всем — и по этой причине вам придется научиться говорить «нет».

Первый вопрос, который нужно задать себе, когда вас просят помочь: «Хочу ли я работать над этим проектом?» Некоторые проекты могут быть интересны, а другими могут руководить ваши друзья, но в конечном счете ваше решение должно основываться на том, какую пользу для компании вы извлечете*. Затем вы можете сравнить все предложения и выбрать то, которое может наиболее благотворно повлиять на вашу компанию. Всегда помните, что основная работа аналитика данных заключается в получении пользы для компании. Этот определяющий критерий успеха также послужит вам подспорьем, когда придется объяснять другим, почему вы не принимаете участие в их проекте.

Как вы *можете* сказать «нет»? Во-первых, воздержитесь от немедленного ответа. Скажите человеку, от которого поступило предложение, что вы подумаете и сообщите о своем решении в течение конкретного периода времени (я даю себе два рабочих дня).

Если вы хотите пойти более дипломатичным путем, можно проделать один трюк, чтобы убедиться, что вас не считают просто «парнем или девушкой по данным», легко щелкающим числа. Если вы позиционируете себя как *советника*, то, даже если вы не участвуете в проекте, вы можете помочь коллегам с использованием данных. Я считаю, что это намного лучше, чем просто сказать «нет»: вы укажете соотаварищам, над чем поработать. Потратьте всего полчаса на исследование того, что им нужно, — возможно, вы найдете инструмент, который подтолкнет их в правильном направлении.

В итоге они будут воспринимать вас не как жонглера данными, а как консультанта — человека, способного дать ценный совет и помочь в достижении цели. Это принесет пользу еще и вам, поскольку позволит глубже познакомиться с деятельностью компании.

* Такого рода польза необязательно имеет денежное выражение, но может предполагать рост количества подписок клиентов, повышение эффективности и т.д. Главное здесь — понимать, что полезно для вашей компании, на этом должен быть основан ваш ответ.

Вперед!

Сделав восемь шагов, описанных в этой главе, вы не только защитите себя от наиболее распространенных проблем, которые могут возникнуть при реализации проектов в области науки о данных, но и начнете самоутверждаться в компании. Наслаждайтесь этим исследовательским и коммуникативным этапом процесса — на следующей стадии вы окажетесь почти полностью прикованы к своему рабочему столу.

Большинство из нас бывали за границей в странах, на языке которых мы не говорим. Когда основной способ коммуникации недоступен, общение оказывается чрезвычайно сложной задачей. Даже если мы немного владеем языком, слабое знание лексики и грамматики часто мешает нам (и слушателю) хорошо понимать друг друга.

Язык, таким образом, фундаментальная необходимость в случае, когда мы хотим понять и общаться с другим человеком. А подготовка данных все равно что создание общего для человека и машины языка.

В этой главе мы узнаем, почему данные никогда не должны анализироваться без предварительной подготовки, каков пошаговый процесс подготовки данных и в чем состоят лучшие методы решения проблем, связанных с массивами данных.

Как заставить данные говорить

К нам как к практикам, если только мы не супервезунчики, данные часто будут попадать «грязными». Нередко данные собираются сотрудниками, не стандартизирующими свои записи, или управляются людьми, которые могут изменить названия столбцов и строк массивов данных в соответствии со своими собственными проектами. Они могут храниться в неподходящих местах — там, где есть риск повреждения. Неудивительно, что при таком множестве разных людей, работающих с одним массивом данных и добавляющих их различными методами, итоговые массивы данных во многих организациях полны ошибок и пробелов. И мы не можем ожидать, что машина будет знать, где находятся ошибки или как исправить несоответствия в информации.

Поэтому наша задача — подготовить данные таким образом, чтобы они были поняты и правильно проанализированы машиной.

С большими возможностями приходит большая ответственность

Подготовка данных (или преобразование сырых данных) является сложным компонентом всего процесса, поскольку она включает в себя ряд задач, которые могут быть выполнены только вручную. Этот этап обычно занимает наибольшее количество времени*. Причина такого пристального внимания к подготовке данных заключается в том, что если исходные данные в массиве изначально не структурированы должным образом, то на более поздних этапах процесс либо вообще остановится, либо, что еще хуже, мы получим неточные прогнозы и/или неправильные результаты. Это может означать катастрофу для вас и вашей компании, и в самом худшем варианте пренебрежение данным этапом может привести к увольнению, а в случае привлечения фрилансеров — даже к судебным искам.

Я не собираюсь пугать вас, просто хочу показать, насколько важно подготовить данные. Удивительно, но, несмотря на важность этого шага, я обнаружил, что учебные материалы науки о данных в основном сосредоточены на более поздних этапах процесса: анализе и визуализации. В этих книгах и курсах используются уже подготовленные массивы данных. Но такой подход хорош, только если вы просто знакомитесь с дисциплиной; в противном случае он означает, что вы эффективно изучаете лишь косметические способы работы с данными.

Работая исключительно с массивами данных из образовательных курсов, вы просто увидите данные, уже очищенные так, как того требует рассматриваемый пример. Но в реальном мире данные часто грязные, перепутанные и поврежденные, и, не зная причин и характеристик грязных данных, мы не можем надлежащим образом

* Мнения аналитиков данных расходятся, но большинство считает, что на подготовку данных уходит 60–80% времени, потраченного на реализацию всего проекта.

завершить проект. Если вы не подготовите данные, то, когда выйдете в реальный мир со своим первым проектом, ваш алгоритм неизбежно выдаст ошибки «отсутствия данных», или ошибки «текстового спецификатора», или «деление на ноль», и проект застопорится.

Но как тогда понять, что данные хорошо подготовлены? Легко, нужно лишь убедиться, что они подходят для нашей стадии анализа данных. Они должны:

- быть правильно отформатированы;
- не иметь ошибок;
- учитывать все пробелы и аномалии.

Распространенная фраза, которую используют аналитики данных, «мусор внутрь, мусор наружу» означает, что если вы примените алгоритм к грязным данным, то получите только бессмысленные результаты, делающие ваш анализ бесполезным. Правда и то, что некоторым практикам с трудом дается этот этап, но только потому, что у них нет шаблона, которому надо следовать. В итоге такие специалисты работают бесструктурно и вынуждены изобретать велосипед каждый раз, когда готовят данные; в долгосрочной перспективе это неэффективный и затратный по времени подход.

Итак, приступим к процессу подготовки данных.

Кейс: Ubisoft — обоснование необходимости подготовки данных

Ульф Морис — финансовый директор немецкого филиала Ubisoft, компании по дизайну, разработке и распространению игр, создавшей популярные игровые франшизы от Assassin's Creed до Far Cry. Ульф курирует дистрибьюторскую дочернюю компанию, продающую видеоигры Ubisoft в Германии, Швейцарии и Австрии (GSA), а также отвечает за финансовые аспекты деятельности компании в Центральной Европе.

Раньше данные в Ubisoft использовались исключительно ее производственной командой для монетизации и внутриигровой аналитики. До тех пор, пока Ульф не изменил ситуацию, финансы не входили в число стратегически важных

областей науки о данных*. Но игнорирование пользы науки о данных может оказаться дорогостоящим просчетом, и Ульф, имевший опыт применения данных при принятии важных бизнес-решений (на предыдущей работе он сохранил компании \$40 млн в слиянии благодаря своему вниманию к данным), знал, что продуманная стратегия использования компанией данных чрезвычайно важна.

Вот что он заявляет:

«Подготовка данных не добавляет вам данных, она просто улучшает способы их исследования. Это похоже на сцену из фильма “Волшебник страны Оз”: Дороти открывает дверь своего дома, попадает в королевство Оз — и черно-белый мир Канзаса превращается в цветной. Вроде бы мало что изменилось на техническом уровне, но все же картина стала другой. Мир приведен в порядок».

(SuperDataScience, 2016)

Чтобы лучше понять, что Ubisoft выиграла от подготовки данных, Ульф обратился к производственному отделу компании. В течение многих лет производственники собирали данные о тысячах онлайн-геймерах, о продолжительности игры и времени, которое требовалось для прохождения отдельных уровней, о том, что игрок делал в игре и в какой момент игры терпел неудачу. Ульф обнаружил, что они используют эти прошлые данные для оценки вероятности того, что клиенты будут покупать внутриигровые предметы через «модель freemium»**. Наличие под рукой данных Ubisoft не только помогло выяснить характер покупок основных клиентов, но и определить поведение, которое можно было ожидать от будущих игроков.

Результаты, к которым производственная команда пришла благодаря науке о данных, заставили всех повернуться лицом к этой идее. На собрании по финансовой стратегии команда Ульфа наметила источники всех доступных данных Ubisoft и то, чего не хватало для полноты картины, и предложила нечто осязаемое, от чего коллеги могли «отщипнуть» во благо своих идей.

«Все очень просто, — говорит Ульф. — Если вы не знаете о существовании чего-то, вы не можете задавать вопросы об этом».

(SuperDataScience, 2016)

* Нередко крупные организации, собирающие данные в течение многих лет, страдают от институциональной слепоты по отношению к науке о данных, не зная, что данные должны быть подготовлены до того, как их можно будет проанализировать, — иначе их информация непригодна для использования.

** Игра загружается бесплатно, но игровые предметы для игроков, которые хотят продвигаться в игре быстрее, продаются за деньги.

Пробелы показали, какие ключевые данные им нужно собрать от своих клиентов (размеры магазина, пространство, отведенное для продажи видеоигр, отношение типичных покупателей к видеоиграм и отношение потребителей к продукции Ubisoft), прежде чем они смогут провести содержательный анализ. Ульф говорит:

«Было очевидно, почему мы не применяем более системный подход к нашим клиентам: у нас нет данных. Отдел продаж оценил для меня клиентов, основываясь на сведениях, не собиравшихся систематически. Получение данных — неопровержимых фактов — было абсолютно необходимо».

(SuperDataScience, 2016)

Сбор этой информации от 2000 магазинов позволил Ульфу подготовить статистически значимые данные, которые в конечном итоге оказались пригодны для анализа. Это помогло Ubisoft выявлять целевую аудиторию как никогда эффективно.

Подготовка данных к путешествию

Для того чтобы сделать исходные (сырые) данные пригодными для анализа, их нужно сначала подготовить:

- 1) извлечь** данные из исходных источников;
- 2) перевести** данные на понятный язык, чтобы они стали доступны в реляционной базе данных;
- 3) загрузить** данные в конечный источник.

Этот процесс известен как ETL (Extract — Transform — Load), и он поможет собрать данные подходящего формата в конечном источнике («хранилище»), к которому можно получить доступ и проанализировать данные на более поздних этапах процесса их обработки. Хранилище содержит разрозненные данные в одной системе. Зачастую оно будет включать реляционные базы данных.

Что такое реляционная база данных?

Реляционные базы данных (РБД) позволяют исследовать их реляционные данные. В таких базах данных имеют значение отношения между единицами информации во всем массиве данных.

Массивы данных в РБД связаны столбцами с одинаковыми именами. Например, если несколько массивов данных содержат столбцы с наименованием «страна», данные из этих столбцов можно сравнить в реляционной базе данных. Преимущество такой базы данных в том, что в ней больше возможностей для анализа и визуализации, необходимых для получения полезных выводов. В частности, данные в такой базе могут изучаться в нескольких массивах сразу без необходимости индивидуального извлечения.

Возможно, лучший способ проиллюстрировать преимущества реляционной базы данных — сравнить ее с Excel, которая часто используется теми, кто не привык работать с базами данных:

- 1. РБД поддерживает целостность.** Каждая ячейка в Excel индивидуальна; типы значений, которые можно в нее поместить, не ограничиваются. Вы можете добавить даты или текст, например, под номерами телефонов или денежными величинами, и Excel это будет полностью устраивать. А вот реляционная база данных станет бить вас по рукам за такую небрежность. Типы столбцов в базе данных предопределены, что означает, что столбец, настроенный на запись дат, не будет принимать значения, не отвечающие формату даты. Затем базы данных будут следить за процессом, делая запрос по любому показателю, который не соответствует значению, предопределенному столбцом.
- 2. РБД комбинирует массивы данных.** Объединить массивы данных в реляционной базе данных легко; гораздо труднее это сделать в Excel. Реляционные базы данных были разработаны для этой цели, и они позволяют легко создавать новые массивы данных путем объединения общих значений в РБД. Все, что от вас требуется, — это умение выполнить простую команду. Поскольку комбинирование таблиц не является основной функцией Excel*, для объединения данных в одну таблицу там требуются расширенные навыки программирования.

* Для различных таблиц используются вкладки, но объединить значения через них может быть сложно.

3. РБД масштабируема. Реляционные базы данных были специально разработаны для масштабируемости; поскольку они объединяют массивы данных, ожидается, что они должны быть в состоянии справиться с большим количеством информационных единиц. Что означает — независимо от того, есть ли у вас пять или пять миллиардов строк, — ваша реляционная база данных вряд ли рухнет в критический момент. Excel гораздо более ограничена в плане емкости, и по мере роста массива данных производительность программы ухудшается, поскольку она изо всех сил пытается справиться с перегрузкой.

Очистка данных

Мы знаем, что в реальном мире данные, скорее всего, будут поступать к нам грязными, но среди практиков есть некоторые разногласия относительно того, как и когда их нужно очищать. Одни очищают данные перед их преобразованием, а другие — только после загрузки в новую базу данных. Я предпочитаю очищать данные на *каждом этапе* процесса ETL — это может показаться неэффективной тратой времени, но я обнаружил, что нет лучшего способа защититься от неприятностей в дальнейшем. К сожалению, подготовка данных всегда будет занимать много времени, но чем больше осмотрительности вы проявите на этом этапе, тем больше ускорите процесс анализа данных в целом.

1. Извлеките данные

Нам нужно извлечь данные: 1) чтобы убедиться, что мы не изменяем каким-либо образом исходный источник; и 2) потому что данные, которые мы хотим проанализировать, часто хранятся в разных местах. Некоторые примеры возможных местоположений:

- база данных;
- таблицы Excel;

- сайт;
- Twitter;
- CSV-файл;
- бумажный отчет.

Если мы используем данные из нескольких источников, нам придется извлечь их в единую базу данных или хранилище, чтобы проанализировать. Но их не всегда легко извлечь из мест, которые используют форматирование, специфическое для конкретной системы, — например, из Excel, к которой мы вернемся позже в этой главе.

CSV-файлы

Как специалист по данным, вы познакомитесь с этими типами файлов довольно близко. Это самый простой тип необработанных файлов с данными, полностью лишенными какого-либо форматирования, что делает их доступными для любого количества программ, в которые мы можем их импортировать. В CSV-файлах строки размещаются на новых строках и столбцы разделяются запятыми в каждой строке. Отсюда и аббревиатура, которая расшифровывается как **comma separated values** (данные, разделенные запятой).

Прелесть работы с необработанными файлами заключается в том, что вы никогда не потеряете или не повредите информацию при загрузке массива данных в программу. Именно поэтому они являются стандартом для большинства практиков.

Почему важно извлекать данные, даже если они находятся только в одном месте

Технически вы можете анализировать данные непосредственно в пределах их хранилища (исходная база данных, электронная таблица Excel и т. д.). Хотя этот метод не рекомендуется, он приемлем для быстрых вычислений, таких как вычисление суммы столбца значений в Excel. Тем не менее для серьезных проектов в области науки о данных работать с данными в их первоначальном хранилище запрещено.

Иначе вы можете случайно изменить необработанные данные, что поставит под угрозу вашу работу.

И это *наилучший* сценарий, поскольку он затрагивает только вас и ваш индивидуальный проект. Работа в хранилище вместо извлечения исходных данных в тестовую базу делает данные уязвимыми для повреждения пользователями, и ваша работа может даже привести к сбою внутренних систем учреждения. Необходимо взять паузу, прежде чем начать работать с данными организации. Нам, аналитикам данных, доверяют важную, существенную информацию о компании, поэтому мы должны убедиться, что оставляем данные такими же, какими они были, когда мы приступили к проекту.

Программное обеспечение для извлечения данных

Для извлечения и чтения данных существует несколько бесплатных программ, и они обязательно отучат вас от вредных привычек, которые часто формируются у пользователей Excel. Эти программы хорошо работают с данными, которые находятся в необработанном файле формата CSV*.

Хотя это может занять некоторое время, данные в большинстве случаев могут быть урезаны до необработанных CSV-файлов. И если вы работаете в большой организации, где вам нужно подать запрос на извлечение данных, то вот хорошие новости: данные, скорее всего, в любом случае будут предоставлены вам в формате CSV.

Notepad++ — инструмент, которым я пользуюсь, когда хочу посмотреть извлеченные мной данные. Это мощный редактор для просмотра CSV-файлов, и он гораздо удобнее, чем программа «Блокнот», которая стандартно поставляется с Windows. Notepad++ также имеет несколько других существенных преимуществ, таких как:

- нумерация строк, позволяющая перемещаться по файлам и отслеживать вкладки с возможными ошибками;

* По мере развития вашей карьеры в науке о данных вы научитесь работать с различными хранилищами данных. Здесь мы говорим о CSV-файлах, потому что они наиболее распространены и универсальны и с них удобно начинать.

- функция поиска и замены, дающая возможность быстро находить значения или текст, которые не нужны в массиве данных, и изменять их;
- Notepad++ был разработан специально, чтобы вы были уверены, что ваши данные не могут случайно измениться, как это может случиться в других программах электронных таблиц;
- в то время как текстовый редактор «Блокнот», поставляющийся с Windows, как правило, имеет проблемы с большими файлами, Notepad++ может открывать файлы размером до 2 ГБ.

EditPad Lite — бесплатная программа для личного использования. Она предлагает возможности, аналогичные Notepad++, но с одним важным преимуществом: хотя обе они хорошо работают с файлами размером до 2 ГБ, я заметил, что Notepad++ иногда может «сопротивляться» массивам данных, близким к максимальному размеру файла. В результате я обнаружил, что EditPad Lite работает с большими файлами намного лучше. Если вы заметите, что перегрузили файлами Notepad++, обратитесь к EditPad Lite.

2. Преобразуйте ваши данные

Нельзя просто сбросить данные из исходного источника непосредственно в хранилище данных — если только вы не *хотите* работать с беспорядочным массивом данных. Преобразовав данные, можно «перевести» информацию, которую планируется использовать, на язык, соответствующий поставленным целям.

В широком смысле этап преобразования включает такие изменения, как объединение, разделение и агрегирование данных. Эти функции позволяют создавать производные таблицы, лучше согласующиеся с имеющейся задачей. Но самая важная функция преобразования — очистка данных, и именно на ней мы сосредоточимся.

На этом этапе мы должны выявить и устранить в нашей исходной базе данных любые ошибки и изъяны, которые часто охватывают весь

спектр — от несоответствий форматирования и резко отклоняющихся значений до значительных пробелов в информации. Но чтобы сделать это, мы сначала должны понять, что мы ищем. Итак, как мы можем выявить грязные данные?

Грязные данные

Грязные данные — это неверная, поврежденная или отсутствующая информация.

Неверные данные — результат того, что информация была (частично или полностью) неправильно добавлена в базу данных (например, ввод значения валюты в ячейку даты). Иногда мы видим, что данные неверны. Это может быть очевидно при несоответствии между столбцами.

Например, если бы у нас была одна строка, где в ячейке страны значилась «Франция», а в ячейке города — «Рим», мы бы поняли, что она неверна. Мы также можем определить неправильные данные, ориентируясь на здравый смысл: так, мы бы знали, что запись в столбце даты рождения в виде «12/41/2001» просто не может быть правильной.

Поврежденные данные — информация, которая изначально в массиве данных была правильной, но оказалась искажена. К факторам порчи информации относятся физическое повреждение базы данных, ее изменение другим программным обеспечением или предшествующее извлечение данных нерекомендуемыми способами. Иногда данные могут просто быть повреждены из-за переноса в базу данных, не поддерживающую формат, который они имели в предыдущем хранилище.

Пропущенные данные возникают, если для данной ячейки нет доступной информации или если лицо, ответственное за вставку данных, не добавило их в ячейку. Пропущенные данные — частое явление в науке о данных, и, вероятнее всего, оно обусловлено человеческим фактором.

Что может произойти, если мы не восстановим недостающие данные

Мы всегда должны знать о любых пробелах в нашей информации. Ниже вы увидите реальный пример данных, извлеченных из электронной таблицы Excel в CSV-файл, который показывает (по годам) выплаты дивидендов (рис. 5.1).

```

487 19-May-15,533.98,540.66,533.04,537.36,537.36,1966900
488 18-May-15,532.01,534.82,528.85,532.3,532.3,2003400
489 15-May-15,539.18,539.27,530.38,533.85,533.85,1971300
490 14-May-15,533.77,539,532.41,538.4,538.4,1403900
491 13-May-15,530.56,534.32,528.66,529.62,529.62,1252300
492 12-May-15,531.6,533.21,525.26,529.04,529.04,1634200
493 11-May-15,538.37,541.98,535.4,535.7,535.7,905300
494 08-May-15,536.65,541.15,525,538.22,538.22,1527600
495 07-May-15,523.99,533.46,521.75,530.7,530.7,1546300
496 06-May-15,531.24,532.38,521.09,524.22,524.22,1567000
497 05-May-15,538.21,539.74,530.39,530.8,530.8,1383100
498 04-May-15,538.53,544.07,535.06,540.78,540.78,1308000
499 01-May-15,538.43,539.54,532.1,537.9,537.9,1768200
500 30-Apr-15,547.87,548.59,535.05,537.34,537.34,2082200
501 29-Apr-15,550.47,553.68,546.91,549.08,549.08,1698800
502 28-Apr-15,554.64,556.02,550.37,553.68,553.68,1491000
503 27-Apr-15,563.39,565.95,553.2,555.37,555.37,2398000
504 26-Apr-15,10000000/10000000 Stock Split,,,.
505 24-Apr-15,564.55,569.58,555.72,563.51,563.51,4932500
506 23-Apr-15,539.52,549.45,538.75,545.5,545.5,4184800
507 22-Apr-15,532.94,539.6,530.29,537.89,537.89,1593500
508 21-Apr-15,536.04,537.91,532.21,532.51,532.51,1844700
509 20-Apr-15,524.16,534.62,523.06,533.91,533.91,1679200
510 17-Apr-15,527.21,528.39,519.58,522.62,522.62,2151800
511 16-Apr-15,528.45,534.12,528.16,532.34,532.34,1299800
512 15-Apr-15,527.25,533.27,521.79,531.07,531.07,2318800
513 14-Apr-15,534.78,536.1,526.65,528.94,528.94,2604100

```

Рис. 5.1. Импорт таблиц, в которых пропущены данные

Как вы можете видеть, часть запятых не разделяют никакую информацию, то есть в пяти столбцах в выделенной строке 504 (26-Apr-15) отсутствуют поля данных.

Нам повезло, что отсутствующие столбцы пережили извлечение — часто отсутствующие значения данных не обрамляются запятыми. В таком случае при использовании на массиве данных алгоритма данные были бы откалиброваны неправильно, в результате чего данные в строке ниже оказались бы сдвинуты для того, чтобы соответствовать требуемому количеству столбцов массива данных. Здесь это означало бы, что дата 24-Apr-15 будет выведена в столбец непосредственно справа от значения «10000000/10000000 Stock Split».

Подобное отсутствие данных может вызвать значительные проблемы на этапе анализа, если мы не отловим проблему заранее. Я знал некоторых неопытных аналитиков данных, которые проверяли верхние 100 строк своего массива данных, но это ошибка новичка: если есть ошибки, вы с гораздо большей вероятностью увидите их в конце массива данных, потому что упущения будут сдвигать информацию.

Исправление поврежденных данных

Чтобы исправить поврежденные данные и сделать их доступными для прочтения машиной, мы можем сначала попробовать следующее:

- повторно извлечь их из исходного файла, чтобы увидеть, не был ли файл поврежден во время первого извлечения;
- поговорить с сотрудником, ответственным за данные, чтобы узнать, может ли он пролить свет на то, как должны выглядеть эти данные, или
- исключить из анализа строки, содержащие поврежденные данные*.

Общение с коллегами

Если вы оказались в ситуации, когда вам не хватает данных и необходимо повторить ваши шаги, чтобы получить дополнительные исходные данные и тем самым обеспечить продвижение проекта, советую поступать следующим образом:

* Ваше решение в конечном счете будет зависеть от того, нужны ли вам данные, и на этот вопрос можно легко ответить, если вы нашли время, чтобы определить вопрос на первом этапе процесса обработки и анализа данных.

- Всегда будьте вежливы с сотрудниками, которые дают вам данные. Некоторые могут быть раздражены вашими усилиями по сбору данных, что будет проявиться в их манере общаться, но вы должны постараться оставаться нейтральным. Помните, что эти люди не являются аналитиками данных и не могут испытывать такую же радость в процессе их сбора, как вы! Объясните им, что проекты, основанные на данных, имеют разные результаты и требуют разных типов данных. Возможно, вам придется несколько раз обратиться к команде, ответственной за массивы данных, поэтому будьте дружелюбны и сделайте их своими единомышленниками.
- Убедитесь, что все, с кем вы общаетесь, полностью понимают проблему, которую вы пытаетесь решить, а также свою роль в этом процессе. Видение более широкой картины поможет коллегам с большей терпимостью относиться к вашим запросам.
- Всегда имейте под рукой список информационных активов компании. Когда вы отправитесь на охоту за новыми данными, он пригодится, чтобы выявить то, что у организации уже есть, и уменьшит вероятность повторного сбора одних и тех же данных. Я рекомендую в этом перечне записывать названия источников, а также столбцов баз данных и их дескрипторов.

Восполнение недостающих данных

Если мы не можем решить проблему, используя любой из этих методов, то придется рассматривать часть данных как отсутствующие. Существуют различные способы решения проблемы пропущенных данных в электронных таблицах:

- **Точно определите, какая именно информация отсутствует.** Это можно сделать для информации, полученной из других данных. Например, предположим, что у нас есть электронная таблица с данными о местоположении клиента, которая содержит значения столбцов как для «штата», так и для «города»; запись, соответствующая «штату», отсутствует, но значение

«города» — «Солт-Лейк-Сити». Тогда мы можем быть уверены, что штат — «Юта»*. Также можно получить пропущенное значение на основе нескольких значений, например для получения значения прибыли из разницы доходов и расходов. Имейте в виду, что мы вводим информацию в обоих примерах, исходя из предположения, что при сборе данных не было ошибок.

- **Оставьте запись как есть.** Можно просто оставить ячейку без данных незаполненной. Это особенно полезно, если определенные поля не имеют никакого отношения к нашему анализу и, следовательно, могут быть исключены из тестирования. Прием может также использоваться, если мы планируем применить метод, который незначительно пострадает от потери данных (то есть метод, использующий усредненные значения), или если мы используем программный комплекс, который может должным образом преодолеть отсутствие информации. В случаях, когда вы оставляете запись как есть, я бы рекомендовал отмечать, где ваши данные содержат пробелы, чтобы можно было учесть любые возникшие впоследствии аномалии.
- **Полностью удалите запись.** Иногда недостающие данные имеют решающее значение для анализа. В этом случае подходит один-единственный способ — удаление из анализа всей строки, так как недостающая информация делает данные непригодными для использования. Очевидно, однако, что результаты станут менее значимыми по мере уменьшения выборки. Таким образом, этот подход, вероятно, лучше всего работает с большими массивами данных, где пропуск одной строки не сильно повлияет на статистическую значимость всего массива данных.
- **Замените отсутствующие данные средним/медианным значением.** Это популярный подход для столбцов, содержащих числовую информацию, так как он позволяет произвольно восполнять любые пробелы, не внося значительных изменений в массив данных. Чтобы вычислить среднее, мы складываем

* Будьте осторожны с подобными полями. В Соединенных Штатах есть только один Солт-Лейк-Сити, но иногда вы найдете несколько городов, называющихся одинаково.

все значения и делим сумму на количество значений. Чтобы вычислить медиану, мы находим последовательное среднее значение в нашем диапазоне данных (если число значений нечетное, просто сложите два средних числа и разделите сумму на два). Обычно предпочтительнее вычислять медиану, а не среднее значение, поскольку первая меньше подвержена влиянию резко отличающихся значений, а это означает, что экстремальные значения по обе стороны от медианного диапазона не будут искажать результаты.

- **Заполните пропуски, исследуя корреляции и сходства.** Этот подход снова зависит от числового значения отсутствующих данных и требует использования моделей прогнозирования возможных пропущенных значений. Например, мы могли бы использовать прогностический алгоритм (скажем, алгоритм k-ближайших соседей, который мы обсудим в главе 6) для вставки недостающих данных на основе существующих сходств между записями в нашем массиве данных.
- **Введите фиктивную переменную для отсутствующих данных.** Это требует добавления столбца в наш массив данных: везде, где мы находим пропущенные данные, мы присваиваем ячейке значение «да» — а когда они не пропущены, даем ей значение «нет». Затем мы можем изучить, как переменная коррелирует с другими значениями в нашем анализе, и ретроспективно рассмотреть возможные причины отсутствия этих данных.

Действия в случае наличия резко отклоняющихся значений

Предположим, что мы работаем на компанию, продающую аксессуары для телефонов, и хотим найти среднее количество чехлов одной модели, проданных каждому из наших дистрибьюторов. Мы работаем уже много лет, поэтому у нас большие массивы данных. У сотрудника, ответственного за ввод этих значений в базу данных, был плохой день, и, вместо того чтобы ввести в столбец «продукт» количество единиц продукта, он вставил туда номер телефона дистрибьютора. Эта ошибка аномально увеличила наш средний показатель

в этой колонке (и означала бы, что один дистрибьютор купил по крайней мере 100 млн единиц продукта!). Если бы мы проанализировали эту запись отдельно, то, вероятно, заметили бы ошибку. Но если бы мы просто рассчитали среднее значение, не глядя на данные, наш отчет был бы искажен этим резко отклоняющимся значением — и это сделало бы его непригодным.

Тем не менее важно различать резко отклоняющиеся значения, которые могут быть отнесены к ошибочной информации, и те, что являются правильными, но выходят за пределы нормального диапазона значений. Если дистрибьютор действительно купил 100 млн единиц вашего продукта, это значение все равно будет резко отклоняться, поскольку оно выше стандартного.

Многие массивы данных содержат резко отклоняющиеся значения, и наша задача — понять, где они находятся, и убедиться, что они не искажают фатально наши отчеты. Это во многом будет зависеть от того, какой анализ мы хотим провести. Например, если бы мы хотели выяснить для издательства среднее количество единиц, проданных книжным магазинам по всему миру, и при этом знали, что резко отклоняющееся значение связано с исключительным заказом на поставку, мы могли бы удалить запись, даже если она достоверна.

Можно найти резко отклоняющиеся значения в массиве данных без их поиска вручную путем создания кривой распределения (также известной как колоколообразная кривая нормального распределения) на основе значений столбцов. Кривые распределения графически отображают на пути к их вершине наиболее вероятное значение или событие из ваших данных, и их достаточно просто создать, даже в Excel*. После создания кривой распределения можно определить значения, выходящие за пределы нормального диапазона.

Кейс: прикладные методы работы с грязными данными

В предоставленном нам массиве данных из воображаемого фонда венчурного капитала (рис. 5.2) отражался общий рост стартапов в Соединенных Штатах.

* Версии MS Office различаются. Ввод «кривая распределения» в меню справки Excel покажет результаты, необходимые для создания кривой.

Поскольку сборщик данных не был связан со стартапами, некоторая информация отсутствовала, так как она либо не была общедоступной, либо компании-стартапы не желали предоставлять информацию такого уровня.

ID	Name	Industry	Inception	Employees	State	City	Revenue	Expenses	Profit	Growth
1	Over-Hez	Software	2008	25	TN	Franklin	\$9,684,527	1,130,700 Dollars	8553827	19%
2	Unimattax	IT Services	2009	35	MI	Kirklin	\$14,016,543	804,035 Dollars	13212508	20%
3	Greentax	Retail	2012	35	MI	Procy with industry Median	\$9,746,272	1,044,375 Dollars	8701997	16%
4	Blacklane	IT Services	2011	64	CA		\$15,559,369	4,631,808 Dollars	10727561	19%
5	Yearflex	Software	2013	45	CA		\$8,567,910	4,374,841 Dollars	4193069	19%
6	Indigoplanet	IT Services	2013	60	CA		\$12,605,452	4,626,275 Dollars	8179177	22%
7	Tresiam	Financial Services	2009	116	MO	Clayton	\$5,387,469	2,127,984 Dollars	3259485	17%
8	Radrimind	Construction	2013	73	NY	Woodside				
9	Lamton	IT Services	2009	55	CA	San Ramon	\$11,757,018	6,482,455 Dollars	5274553	30%
10	Striplnd	Financial Services	2010	25	FL	Boca Raton	\$12,329,371	916,455 Dollars	11412916	20%
11	Canecorporation	Health	2012	6	NY	New York	\$10,597,009	7,591,189 Dollars	3005820	7%
12	Mattouch	IT Services	2013	6	WA	Bellevue	\$14,026,934	7,429,377 Dollars	6597557	28%
13	Techonit	Health	2009	9	MS	Flowood	\$10,573,990	7,435,363 Dollars	3138627	8%
14	Technine		2006	65	CA	San Ramon	\$13,898,119	5,470,303 Dollars	8427816	22%
15	Chyse		2010	25	CO	Louisville	\$9,254,614	6,249,498 Dollars	3005116	6%
16	Kayelelectronics	Health	2009	687	NC	Clayton	\$9,451,943	3,878,113 Dollars	5573830	4%
17	Ganzclax	IT Services	2011	75	NJ	Iselin	\$14,001,180		11901180	18%
18	Trantraolax	Government Services	2011	35	VA	Suffolk	\$11,088,336	5,635,276 Dollars	5453060	7%
19	E-Zm	Retail	2008	320	OH	Monroe	\$10,746,451	4,762,319 Dollars	5984132	13%
20	Darface	Software	2011	78	NC	Durham	\$10,410,626	6,196,409 Dollars	4214219	17%
21	Holiane	Government Services	2012	87	AL	Huntsville	\$7,978,332	5,666,574 Dollars	2291758	2%
22	Lathotline	Health		103	VA	McLean	\$9,418,303	7,567,233 Dollars	1851070	2%
23	Lambam	IT Services	2012	210	SC	Columbia	\$11,950,148	4,365,512 Dollars	7584636	20%
24	Quozap	Software	2004	21	NJ	Collingswood	\$8,304,480	7,019,973 Dollars	1284507	20%
25	Tampware	Construction	2011	13	TX	Houston	\$9,785,982	2,910,756 Dollars	6875286	11%
26	Dalhow	Health	2000	20	GA	Dacula	\$10,600,718	7,731,820 Dollars	3068998	7%
27	Ranktach	Government Services	2010	607	FL	Tampa	\$10,515,557	7,439,384 Dollars	3078173	8%
28	Unadex	Software	2013	280	NC	Cary	\$5,231,275	2,388,521 Dollars	2842754	19%

Рис. 5.2. Отсутствующие данные в таблице для показа роста стартапов

Как вы можете видеть, различные типы информации отсутствуют в столбцах, а иногда в одной строке есть несколько пустых значений. Давайте применим на практике методы исправления недостающих данных. Вернитесь к методам, представленным выше, и подумайте, как бы вы могли решить проблему недостающих данных самостоятельно, прежде чем читать ответы ниже.

Сотрудники

Замените отсутствующие данные средним/медианным значением. Это числовое значение, и поэтому мы можем на место любого из пропущенных значений «сотрудников» использовать общую или отраслевую медиану для этого столбца. (Отраслевой медианный показатель предпочтительнее, поскольку он будет аналогичен отсутствующему показателю.)

Отрасль

Оставьте запись как есть, или точно определите, какая именно информация отсутствует, или полностью удалите запись. Выяснить, к какой отрасли относится компания, можно просто исследуя, что она делает, и на этом построить ваши предположения. Но выбор зависит от того, насколько важна отрасль для нашего анализа. Если отрасль важна, а мы не можем ее определить, нужно удалить запись из анализа.

Год основания компании

Оставьте запись как есть, или точно определите, какая именно информация отсутствует, или полностью удалите запись. Несмотря на то что дата — это число, оно не является числовым значением (с ним нельзя выполнять арифметические операции). Значит, мы не можем заменить его средним значением, а если мы не можем узнать, когда была создана компания, то мы должны воспринимать эту информацию как отсутствующую.

Штат

Оставьте запись как есть, или точно определите, какая именно информация отсутствует, или полностью удалите запись. Мы можем безошибочно предположить, какие сведения должны быть на месте недостающих данных. Но требуется осторожность: в случаях, когда город с таким названием может находиться более чем в одном штате, речь не может идти о 100%-ной точности предложенного значения, и поэтому нам необходимо решить, насколько важны эти данные для нашего анализа.

Расходы

Точно определите, какая именно информация отсутствует. Это легко, мы можем рассчитать расходы, просто вычитая прибыль из дохода.

Доходы, расходы и прибыль, рост

Замените отсутствующие данные средним/медианным значением. Чтобы вычислить эти недостающие данные, требуется больше шагов. Нужно сначала заменить рост доходов и расходов, используя медианы отрасли, а затем мы сможем рассчитать прибыль как разницу между доходами и расходами.

Преобразование данных из MS Excel

Excel пытается упростить задачу, автоматически переформатируя определенные значения. Это может привести к различным сбоям в процессе ETL, и, поскольку программа Excel часто используется для хранения данных, я уделяю ей особое внимание. Одна общая жалоба, которую я слышал от пользователей Excel, — требование программы

преобразовывать длинные числовые значения (такие, как номера телефонов и кредитных карт) в научную формулу*. И это не самое худшее. Excel может конвертировать даты и денежные суммы в единый формат, соответствующий региональным настройкам вашего компьютера. Хотя это может быть удобно для отдельных электронных таблиц, которые часто используются в бизнес-аналитике, такие виды автоматизации в конечном итоге доставят вам неприятности при анализе о данных, так как форматирование Excel не предусматривает качественного перевода в базу данных. И если мы имеем дело с большим количеством данных, выбор всех единиц, измененных программой Excel, может занять много времени.

Если мы не преобразуем данные из Excel в CSV-файл, то в дальнейшем будем сталкиваться с проблемами. Если восстановить измененные даты удастся, то почти невозможно восстановить номера кредитных карт, если они были заменены на числа в экспоненциальной записи. Только представьте, чем это чревато для организации, теряющей номера кредитных карт своих клиентов, особенно если вы работали с единственной копией файла.

Некоторые из наиболее распространенных проблем связаны с датами и валютой, так как их значения не являются международными и поэтому зависят от региональных настроек наших машин.

Форматы дат. Форматирование дат будет отличаться в зависимости от географического региона, и в Excel предустановлено отображение той даты, которая соответствует региональным настройкам нашего компьютера. В большинстве стран используется формат даты, начинающийся со дня, за которым следуют месяц и год (ДД/ММ/ГГГГ). Однако в Соединенных Штатах формат даты начинается с месяца, а потом идут день и год (ММ/ДД/ГГГГ). Необходимо обеспечить согласованный формат дат в нашей базе данных.

Как их исправить. Лучший способ предотвратить внесение изменений в записи программой Excel — поменять все ваши форматы дат на ГГГГ-ММ-ДД, так как это однозначный международный стандарт, не зависящий от региональных правил. В Excel выберите

* Например, 4556919574658621 будет отображаться как 4.55692 E+15.

столбец, который хотите исправить, щелкните его правой кнопкой мыши и выберите пункт «Формат ячеек». В окне «Категории» выберите «Дата». В окне «Тип» вы должны увидеть формат ГГ-ГГ-ММ-ДД. Выберите его и нажмите «ОК». Даты будут изменены.

Форматы валют. Форматы валют также будут зависеть от региональных настроек вашего компьютера. В этих случаях необходимо учитывать не только символ валюты, но и используемые десятичные знаки. Символы валюты должны быть полностью удалены из ваших данных, так как в противном случае они будут читаться как текст. Страны используют различные десятичные знаки для своей валюты: они отделяются либо точкой (например, £30.00 в Великобритании), либо запятой (например, €30,00 в Германии).

Обратите внимание, что это касается как десятичной точки, так и разделителя тысяч. Сумма £30,000 будет читаться как тридцать тысяч фунтов в таких странах, как Австралия, где запятую используют для обозначения тысяч, но ее можно читать как тридцать фунтов в таких странах, как Швеция, где запятую используют для указания десятичных знаков. Базы данных функционируют с системами с десятичной запятой, и любые запятые, включая разделители тысяч, должны быть удалены из данных.

Как их исправить. Мы хотим лишить числа символов и запятых. Если в вашей стране используется система с десятичной запятой, необходимо сначала изменить региональные параметры компьютера, чтобы убедиться, что запятая изменена на точку. Выберите столбец, щелкните его правой кнопкой мыши и кликните «Формат ячеек». В окне «Категории» выберите «Валюта». Снимите флажок «Использовать разделитель 1000», чтобы убедиться, что запятые не будут использоваться; выберите «Нет» из выпадающего списка «Символ» и потом — «2» для количества десятичных знаков. Это удалит лишние символы из наших данных*.

* В русскоязычной версии Excel 2013 нет опции «Использовать разделитель 1000», а вместо категории «Валюта» — формат «Денежный». Порядок действий в русскоязычной версии программы: в окне «Числовые форматы» выберите «Денежный», выберите «Нет» из выпадающего списка «Обозначение» и потом — «2» для числа десятичных знаков. — *Прим. науч. ред.*

3. Загрузите данные

После того как мы преобразовали данные в нужный формат, можно загрузить их в нашу конечную цель: хранилище. Как только этот процесс будет завершен, мы должны вручную просмотреть данные в последний раз, прежде чем пропускать их через машинный алгоритм, чтобы быть абсолютно уверенными, что мы работаем с достаточно подготовленными данными.

Проверка качества после загрузки

Загрузка данных в хранилище иногда может вызывать проблемы. Возможно, вы пропустили очистку некоторых грязных данных на предыдущем этапе или некоторые данные просто были загружены неправильно. По этой причине необходимо научиться перепроверять данные в хранилище.

Ниже приведены приемы проверки качества, которые вы всегда должны применять на этом этапе:

- **Подсчитайте количество строк** в конечном массиве данных и сравните с исходным массивом данных. Если результаты разнятся, вернитесь к исходному массиву, чтобы выяснить, что произошло. К сожалению, иногда самый быстрый способ проверить — просто посмотреть, то есть прокрутить данные строка за строкой. Лучше двигаться снизу вверх, а не сверху вниз, потому что любые ошибки в данных, скорее всего, будут внизу.
- **Проверьте столбцы на асимметричность.** Чтобы полностью обезопасить себя от проблем на этапе анализа, проверьте как верхние 100, так и нижние 100 строк.
- **Проверьте столбцы, подверженные повреждению.** Обычно это относится к датам и балансам — они, как мы установили ранее, наиболее уязвимы.
- **Проверьте текстовые значения.** Если у нас есть текстовые значения в свободной форме, полученные из опросов, в ходе которых респонденты набирали ответы на открытый вопрос, то загрузить такой текст в базу данных может оказаться непростой задачей. Обычно базы данных ограничивают максимальное

количество букв в столбце. Это может привести к отсечению части ответа, в результате чего данные будут отсутствовать, а иногда даже влиять на остальную часть массива данных. Текст свободной формы тоже иногда содержит символы, такие как кавычки, которые базы данных могут не распознать или использовать неправильно, поскольку они являются символами квалификатора.

Подумайте (снова) как консультант

Проверка качества — заключительная часть подготовки данных, так что на этой стадии есть риск сбавить обороты. Проследите, чтобы этого не случилось, ведь обеспечение качества играет важнейшую роль в деле подготовки ваших данных. Мне посчастливилось войти в область науки о данных через мир консалтинга, который уделяет большое внимание проверке качества. На этой стадии работа аналитика данных оценивается коллегами. Цифры должны совпадать, и результаты должны иметь смысл. Не бойтесь этого этапа — он предназначен не для того, чтобы подловить вас, а чтобы помочь защитить вас от ошибок в дальнейшем.

Компании, имеющие опыт работы с данными, разрабатывают строгие процедуры, которым аналитики данных должны следовать буквально, прежде чем проводить какой-либо анализ. В некоторых даже есть консультанты, проверяющие ваши действия столько времени, сколько на это потребуется. Неверный результат будет по меньшей мере стоить денег, а в худшем случае может серьезно повлиять на бизнес-операции. Вот почему так важно убедиться, что контроль качества выполнен, прежде чем перейти к следующему шагу.

Теперь, когда у вас есть прекрасное хранилище кристально чистых данных и вы знаете вопрос или серию вопросов, которые вы хотите им задать, вы можете наконец перейти к моему любимому этапу: анализу.

Анализ данных (часть I)

06

Специалисты, пытающиеся раскрыть принципы науки о данных незнакомым с ней людям, как правило, теряют свою аудиторию, лишь только речь заходит об анализе данных. Но это больше связано с психологией, чем с чем-либо еще: иррациональный страх перед словом «аналитика» автоматически заставляет многих почувствовать себя не в своей тарелке.

То, что я хочу сделать в следующих двух главах, — это не просто объяснить, как работают отдельные алгоритмы, но также выделить *контекст*, в котором их используют аналитики данных. В отличие от справочников, где подробно описываются теоретические тонкости алгоритма, я нахожу, что контекстуальный подход гораздо больше годится для всех типов учеников, независимо от их умений и способностей. В предыдущих двух главах мы рассмотрели: 1) как можно подойти к задаче, касающейся данных, и 2) как подготовить эти данные для анализа. Теперь мы наконец готовы приступить к анализу данных. Этот этап на самом деле описывается подходом «подключи и работай». Самая трудная часть процесса анализа данных состоит в том, чтобы понять многоаспектность проблемы и учитывать ее переменные. После того как мы поняли задачу и вопросы, которые надо поставить, применить алгоритм для ответа на эти вопросы должно быть проще простого.

Не пропустите этот шаг

Даже если вы считаете, что у вас нет необходимых инструментов или математических способностей, не позволяйте себе соблазниться и пропустить следующие две главы, думая, что можете просто нанять

кого-то, чтобы сделать предварительный анализ за вас. Искушенность в математике или другой научной дисциплине может быть полезной в этот момент, но не является обязательным условием. И хотя вы можете добиться успеха, просто зная, как представлять, готовить и собирать данные, все равно нужно по крайней мере понимать каждый этап процесса, чтобы стать профессиональным аналитиком данных.

Самые основные алгоритмы, используемые в анализе данных, которые мы обсудим в главах 6 и 7, можно разделить на три группы*:

- 1) алгоритмы классификации;
- 2) алгоритмы кластеризации;
- 3) алгоритмы обучения с подкреплением.

Используя эти алгоритмы, мы можем понять, как могли бы начать детализировать данные, разрабатывая идеи, которые, возможно, не были очевидны при визуальном анализе. В этой главе мы будем использовать первые две категории: алгоритмы классификации и кластеризации. Хотя важно отметить, что это только две ветви анализа, провести классификацию и кластерный анализ позволяют относительно простые и часто используемые алгоритмы, которые помогут вам быстро работать с данными.

Информация vs математика в науке о данных

Большинство методов, которые мы обсуждаем в этой книге, основаны на сложной математике и статистике. Однако вы могли заметить, что в их описании отсутствуют математические формулы. Это вызывает тревогу: сможем ли мы действительно понять алгоритм, если не вникаем в его детали?

* Я выбрал эти группы потому, что они, как мне кажется, самые важные семейства алгоритмов. Все алгоритмы не самые простые (например, алгоритмы регрессии) и не самые сложные (такие, как нейронные сети и глубокое обучение). На мой взгляд, это самые полезные примеры анализа данных, которые вы можете применить, читая эту книгу.

Вот как я отвечаю на этот вопрос: подумайте о вождении автомобиля. Вам доводилось когда-нибудь разбирать машину? Вы можете отличить распределительные валы от коленчатых? Как на самом деле работает круиз-контроль машины? Большинству из нас неизвестно все, что касается технической эксплуатации наших автомобилей, и все же почти все мы ездим на них. Часто. В этом разница между математикой и интуицией.

Математика разбивает алгоритм на части, чтобы понять, как именно он работает и почему. В этом нет ничего плохого, и бывают ситуации, когда требуется такой уровень детализации. Но по большей части в нем нет необходимости при работе в качестве аналитика данных. Так же как базовое умение пользоваться педалями и рулевым управлением автомобиля поможет вам добраться из пункта А в пункт В, так и интуиция, лежащая в основе аналитических моделей в науке о данных, окажется полезной для решения поставленной задачи.

Если это вас немного успокоило, значит, я выполнил свою работу. Слишком часто я сталкиваюсь с тем, что науку о данных чересчур усложняют. Моя цель — доказать вам, что, как любой человек может водить машину, любой может быть аналитиком данных.

Классификация или кластеризация?

Давайте прежде всего различать эти две категории. Проще говоря, мы используем классификацию, когда уже знаем, в какие группы хотим объединить наши данные с помощью анализа, и мы используем кластеризацию, когда *не знаем*, что это будут за группы с точки зрения чисел или названия. Например, если бы мы хотели провести анализ ответов «да/нет» на вопрос, мы бы использовали алгоритм классификации, потому что знаем, какими будут две результирующие группы: «да» и «нет». А вот если бы мы хотели оценить респондентов одного и того же опроса на основе их возраста и расстояния до ближайшего магазина нашей компании, то использовали кластеризацию, потому что группы результатов, которые будут полезны для нас, не могут быть точно определены заранее (если мы ранее не проводили такой же анализ).

Предположим, авиакомпания обратилась к нам с просьбой узнать, продолжат ли клиенты пользоваться ее услугами или нет (будет ли так называемый «отток клиентов»). Поскольку компания собрала данные об ответах клиентов и их перемещениях (частота полетов, пункт назначения, класс судна, использование услуг на борту, запросы на перевозку багажа), мы можем использовать эти переменные для определения поведения, которое в наибольшей степени свидетельствует о намерении клиента отказаться от услуг авиакомпании. В этом случае мы попытаемся использовать упомянутые факторы, чтобы разделить клиентов на две группы: группа 1 включает тех, кто может прекратить пользоваться услугами авиакомпании, в то время как в группу 2 войдут клиенты, которые, вероятно, продолжат летать самолетами этой авиакомпании. По этой причине мы будем использовать классификацию, потому что распределяем (*классифицируем*) клиентов по двум группам.

И классификация — это то, с чего мы начнем.

Классификация

Если еще *до проведения анализа* мы будем знать, в какие группы попадут наши данные, то тогда лучше пойти по пути классификации, а не кластеризации. В приведенном выше примере клиентов можно рассматривать через записанные о них данные — их обычные маршруты полета, их возможности по тратам, уровень членства в программе для часто летающих клиентов и даже предпочтения мест в салоне. Эти описательные функции могут показаться обширными, но они всего лишь инструменты. Основная цель состоит в том, чтобы классифицировать клиентов так, чтобы они оказались *только в одной из двух групп* — на данный момент компания не заинтересована в поиске чего-либо еще.

При таком анализе важно также иметь предварительные данные, с помощью которых мы можем следить за характеристиками, которые нас интересуют. Это единственный способ создания алгоритма классификации, то есть используя уже имеющиеся примеры.

Следующие алгоритмы классификации расположены в порядке возрастания сложности. Начнем с дерева решений, так как многие читатели уже знакомы со структурными схемами. Такие схемы

используют один и тот же принцип последовательного разделения информации на части, прежде чем представить участнику окончательный ответ. Регрессия по методу случайного леса — это просто расширение алгоритма построения деревьев решений, поскольку в ней используется несколько решающих деревьев для отдельных компонентов массива данных, чтобы обеспечить более точные результаты. Как метод *k*-ближайших соседей, так и наивные байесовские алгоритмы классифицируют точки данных по группам в соответствии с их относительным расстоянием друг от друга, измеряемым переменными каждой записи. Разница между ними станет очевидной в отдельных разделах. Заканчиваем наш обзор классификации логистической регрессией, которая является алгоритмом, используемым, *именно* когда мы хотим оценить вероятность того, что событие произойдет.

Когда вы читаете про эти алгоритмы, имейте в виду мой первоначальный совет учиться, используя *интуицию*: сосредоточьтесь на понимании цели каждого алгоритма и попытайтесь увидеть *предназначение* предпринимаемых шагов. Потратьте время, чтобы переварить каждый из них, — это действительно поможет в данном случае выиграть гонку.

Деревья решений

Дерево решений можно визуализировать в виде блок-схемы. Алгоритм тестирует отдельные атрибуты в массиве данных, чтобы определить возможные результаты, и продолжает добавлять результаты по мере выполнения дальнейших тестов, останавливаясь только тогда, когда все результаты исчерпаны.

Листья этих деревьев дают нам все возможные ответы на все вопросы, которые мы можем задать нашим данным. Мы все порой отвечаем на вопросы журнального теста, когда надо отметить «да» или «нет», чтобы узнать свой тип личности, Леонард вы или Шелдон из «Теории Большого взрыва» или как вы действительно относитесь к йогурту. В этих случаях вопросы — ветви, а результаты — листья.

В мире бизнеса деревья решений можно использовать, скажем, для классификации групп клиентов. Вспомните пример Ubisoft из главы 5: если бы команда разработчиков игр собрала информацию о потенциальном новом подписчике, они могли бы использовать

дерево решений для проверки того, сможет ли он стать участником, на основе массива данных компании о текущих подписчиках. Алгоритм построения дерева решений делит данные об играх компании на листья, которые отражают четкие различия между такими значениями, как время, проведенное за игрой, и возраст, и соотнесет новые данные с одним из результатов, которые мы определили заранее, — в данном случае с «выгодным» для компании и «невыгодным».

Как работают деревья решений

Давайте исследуем эту проблему. Поскольку у нас есть информация о среднем времени игры наших геймеров и их возрасте, мы можем использовать классификацию с помощью дерева решений, чтобы принять относительно их обоснованное решение. Это означает, что в первую очередь нам необходимо иметь следующие данные о текущих подписчиках нашей игровой компании: общее время, проведенное за играми за последний месяц, и возраст*.

Мы создали точечную диаграмму (рис. 6.1) с большим количеством точек данных на основе возраста (X_1) и времени, проведенного за игрой в часах (X_2).

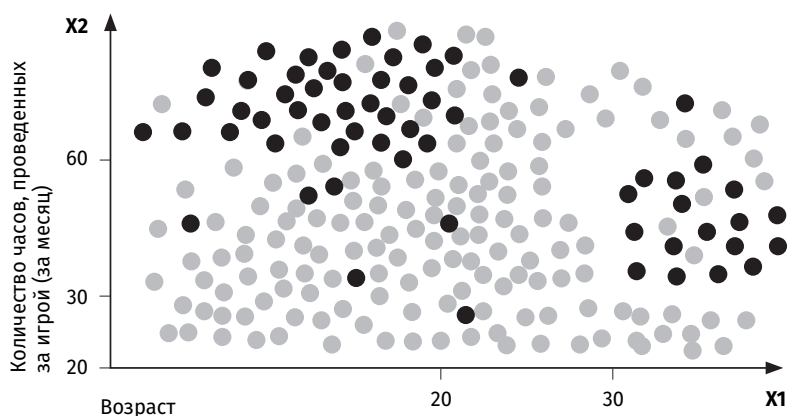


Рис. 6.1. Корреляционный график отображения данных о клиентах, организованных по возрасту и затраченному на игры времени

* Хотя мы будем использовать здесь две переменные, ваш алгоритм не должен ограничиваться только ими.

Серые точки обозначают пользователей, которые не стали подписчиками; черные — подписавшихся. Если бы мы запустили классификационный алгоритм дерева решений, точечная диаграмма была бы разбита на листья, как определено алгоритмом.

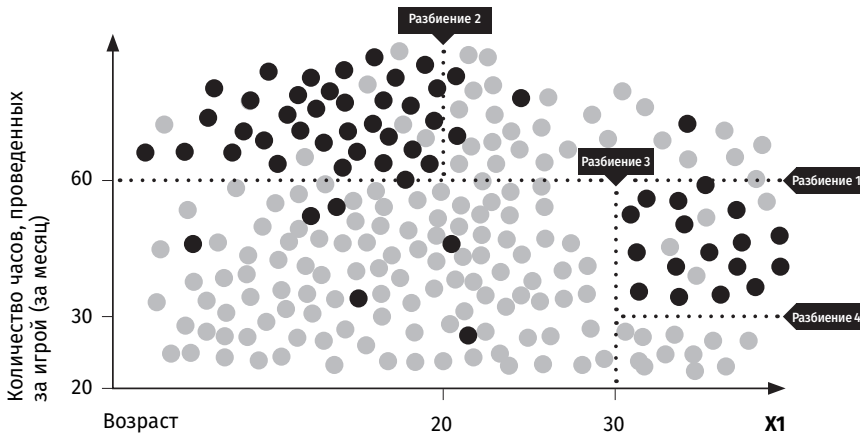


Рис. 6.2. Корреляционный график отображения разделенных на листья данных о клиентах

Как мы видим на рис. 6.2, разбиение 1 разделило данные на $X_2 = 60$, а разбиение 2 — на $X_1 = 20$ и т. д. Алгоритм сгруппировал наши точки данных в листья, что добавляет ценность классификации, и он остановится на оптимальном количестве листьев. Это оптимальное число достигается, когда дальнейшее разбиение данных делает результат листьев статистически незначимым.

Построение алгоритма классификации дерева решений

На рис. 6.2 мы можем проследить логику того, как создается алгоритм дерева решений:

1. Разбиение 1 делит точки данных на те, которые находятся выше и ниже 60 (часов) на оси X_2 .
2. Для тех точек, которые оказываются выше 60 (часов) на оси X_2 , разбиение 2 проводит дальнейшее деление для точек данных на те, которые попадают выше и ниже 20 (лет) на оси X_1 . Это

означает, что разбиение 2 делит только данные, находящиеся выше 60 (часов) на оси X2.

3. Разбиение 3 обращается к точкам данных, которые проигнорировало разбиение 2, разделив те, что оказались ниже 60 (часов) на оси X2. На этот раз разбиение делит точки данных, которые оказываются выше и ниже 30 (лет) на оси X1.
4. Разбиение 4 делит точки данных, находящиеся до 60 (часов) по оси X2 (как разделено с помощью разбиения 1), и те, кому за 30 (лет) по оси X1 (как разделено с помощью разбиения 3). На этот раз разбиение делит точки данных, которые оказываются выше и ниже 20 (часов) на оси X2.

Мы можем перевести описанный процесс в схему на рис. 6.3*.

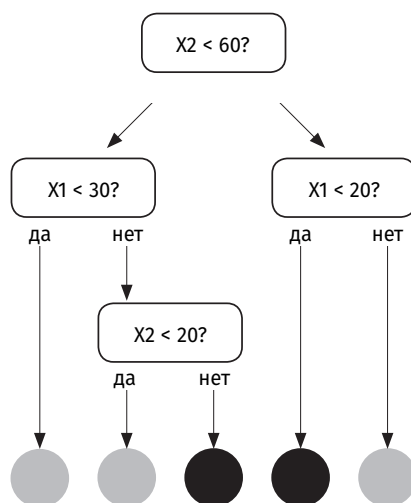


Рис. 6.3. Блок-схема построения дерева решений

Чем этот алгоритм полезен для нас? Предположим, что нашему новому клиенту 19 лет и за первый месяц он потратил 65 часов игрового

* Чтобы избежать глупых ошибок, я считаю, что лучше всего пояснить свой подход последовательной маркировкой. Как вы можете видеть на блок-схеме, все ветви «да» находятся слева, а все ветви «нет» — справа. Все кажется очевидным, но вы будете удивлены количеством людей, которые пренебрегают этим.

времени ($X_1 = 19$ и $X_2 = 65$). На графике рассеяния, который был разделен нашими разбиениями на листья, мы можем видеть, чему будет соответствовать эта точка данных. Наш алгоритм уже обнаружил, что статистически значимое число людей в возрасте до 20 лет, которые играют более 60 часов, с высокой степенью вероятности станут постоянными пользователями, а это значит, что мы можем нацелить на этого клиента рекламу, которая побудит его или ее заплатить за подписку.

Случайный лес

Алгоритм случайного леса основывается на концепции деревьев решений с использованием ансамблевого метода классификации. Вместо одного дерева случайный лес использует много разных деревьев, чтобы сделать один и тот же прогноз, принимая среднее значение результатов отдельных деревьев.

В то время как некоторые из деревьев могут быть недостаточно эффективны (в зависимости от поставленной бизнес-задачи), мы можем применить силу толпы; когда деревья решений используются в совокупности, они могут давать гораздо более обоснованные прогнозы. Подумайте об этом как о процессе голосования — каждое решающее дерево проголосует (сделает прогноз), а затем алгоритм случайного

Способность ансамблевых методов к интуитивному мышлению

Вот мой любимый пример ансамблевого метода, который должен сделать концепцию более интуитивно понятной. Вы когда-нибудь участвовали в конкурсе, в котором нужно определить количество конфет внутри стеклянной банки? Тот, чья догадка окажется ближе всего к правильному ответу, выигрывает приз*.

* Правила игры меняются. Иногда вам нужно угадать точное число конфет, иногда же достаточно дать ответ как можно более близкий к фактическому количеству. Представленная здесь стратегия лучше всего подходит для последнего случая.

Вот моя стратегия того, как сделать более точное предположение на основе ансамблевого метода. Вместо того чтобы просто самостоятельно выбрать число, посмотрите, что отвечают другие участники игры (в большинстве случаев угаданные числа фиксируются на листе бумаги или на доске). Вычислите среднее, округлите результат до ближайшего целого числа, и вуаля! Вот и ваше предположение.

Чем больше догадок вам доступно, тем ближе ваш средний результат будет к фактическому количеству конфет в банке. Это происходит потому, что предположения одних были выше, а других — ниже, но в целом (если сама банка лишена сюрпризов) догадки будут нормально распределены вокруг реального числа конфет. Представьте себе этот эксперимент и поймите, что решение интуитивно. А если вам так не кажется, вы всегда можете попробовать сами.

Если вы замените игру «конфеты в банке» задачей науки о данных и представите предположение каждого из участников деревом решений, то получится ансамблевый метод. Именно таков его принцип: мнение толпы (ансамбля) всегда будет более достоверным — и зачастую более точным, — чем мнение одного респондента (дерево).

леса возьмет вариант с наибольшим количеством голосов в качестве результата. Демократия среди деревьев!

Это делает оригинальный алгоритм намного более мощным. Вместо одного дерева решений для всего массива данных случайный лес создает несколько деревьев решений. Чтобы сделать такие деревья уникальными, их создают из различных подмножеств массива данных.

Давайте исследуем случайный лес чуть подробнее, на примере конкретного случая.

Кейс: BCG — поиск лучших локаций для новых отделений банка

Алгоритмы случайного леса идеально подходят для задач, требующих более комплексной оценки наших данных, чем та, что может быть получена с помощью алгоритма дерева решений. Например, если бы мы хотели оценить потенциал банка,

открывающего филиал в конкретном районе, на основе набора переменных, то мы использовали бы алгоритм случайного леса.

Я живу в Австралии, и, когда хочу зарегистрироваться в новом банке, моим главным приоритетом является удобство. Я хочу, чтобы филиал находился рядом с моим домом, офисом и местом, где я делаю покупки. Если у банка также есть большое количество филиалов рядом с пляжем, еще лучше. Нет ничего хуже, чем ездить на другой конец города, чтобы поговорить с консультантом или обналичить чек.

Банки знают, что удобство — один из основных факторов, влияющих на принятие решений потенциальным клиентом, но они также хотят, чтобы их вновь открытые филиалы были экономически эффективными. Артему Владимирову, ведущему аналитическому консультанту Бостонской консалтинговой группы (BCG), было поручено решить эту проблему для банка — клиента BCG, который хотел развивать свои отделения по всей Австралии.

Сначала Артем проанализировал демографические данные банка, чтобы выяснить количество его клиентов в каждом из районов Австралии. Он увидел, что, поскольку филиалы банка распределены по стране неравномерно, у него нет данных по некоторым австралийским округам. Для того чтобы составить прогнозы по этим местностям, Артему пришлось провести сравнительный анализ данных районов, которые были как «известны», так и «неизвестны» банку, с помощью общедоступной информации о результатах переписи. Применяя такие демографические данные, как средний возраст, гендерная принадлежность, уровень образования и стоимость жизни, Артем смог получить недостающие сведения. Такое использование данных позволило ему рассчитывать на потенциальный успех создания филиалов в новых местах, имевших характеристики, сходные с характеристиками аналогичных районов, которые уже доказали свою выгодность.

Для решения задач банка Артем использовал алгоритм случайного леса:

«Мы взяли всю клиентскую базу из записей данных банка и использовали статистическую модель случайного леса для определения корреляции между рентабельностью клиентов и их демографическими показателями. Прогнозы были сделаны для районов, где у банка уже имелись клиенты, поэтому нам нужно было только перепроверить, будет ли район прибыльным, сопоставив демографические данные».

(SuperDataScience, 2016)

Определив районы, значимые для банка, Артем составил профиль конкурентов компании и количество их филиалов в этих местностях, снова используя случайный лес для определения доли рынка, которую банк занимал по отношению к конкурентам.

Благодаря применению алгоритма случайного леса Артему не нужно было подробно объяснять, какие демографические данные внесли свой вклад в окончательные показатели, это помогло ему обойти проблему защиты персональных данных и показать банку, какие именно области будут наиболее рентабельными для него.

Построение классификации случайного леса

- 1. Выберите количество деревьев, которые хотите создать.** Для многих программ параметр по умолчанию — десять деревьев. Число, которое вы в конечном итоге выберете, будет зависеть от контекста. Меньшее количество деревьев может обусловить менее точные прогнозы. И наоборот, в большинстве случаев можно использовать любое количество деревьев, поэтому нет необходимости беспокоиться о чрезмерно близкой подгонке алгоритма к данным.
- 2. Установите классификатор в тренировочный набор.** Внедрение классификатора случайного леса в тренировочный набор поможет вам в будущем научиться составлять прогнозы для новых точек данных. Затем мы можем сравнить эти прогнозы с фактическими результатами в нашем массиве данных, чтобы увидеть, насколько точен классификатор.

Алгоритм случайного леса случайно выберет N подмножеств из вашего массива данных, где N — количество деревьев, указанное для параметра в шаге 1. Эти подмножества могут перекрываться; однако никакие два множества не будут идентичными.

После выбора подмножеств каждое из них будет использоваться в качестве исходного массива данных для построения уникального дерева классификации. Таким образом, каждое дерево классификации видит только свое подмножество данных и не имеет представления

о том, что фактический массив данных шире. Подобный подход обеспечивает разнообразие при генерации деревьев — именно отсюда в алгоритме случайного леса возникает «сила толпы».

Исходя из этой логики, чтобы помочь алгоритму делать более точные прогнозы, мы можем просто добавить информацию в наш массив данных — чем больше данных в нашем тренировочном наборе, тем более точным будет прогноз алгоритма.

Дерево решений или случайный лес?

Хотя алгоритм случайного леса можно рассматривать как «обновление» деревьев решений, оба метода имеют свои преимущества в зависимости от поставленной задачи. Для проектов, использующих относительно мало данных, применение алгоритма случайного леса не даст оптимальных результатов, так как он будет излишне подразделять данные. В этих сценариях более эффективно дерево решений, которое обеспечивает быструю и простую интерпретацию данных. Но если вы работаете с большим массивом данных, более точный прогноз даст случайный лес, но его интерпретируемость окажется ниже*.

Метод *k*-ближайших соседей (*k*-NN)

Этот метод использует шаблоны в данных для размещения новых точек данных в соответствующих категориях. Предположим, что врач из Сан-Франциско прочитала о недавнем увеличении числа больных диабетом в Соединенных Штатах и хочет учесть эту информацию в своей практике. Врач знает, что сахарный диабет второго типа легче предотвратить, чем лечить, и поэтому она просит аналитика данных разработать на основе записей в медицинских картах ее нынешних пациентов (с диагнозом «диабет» либо здоровых) модель, которая оценит вероятность того, что у ее новых клиентов в будущем может развиться это заболевание. Таким образом, наш врач надеется, что она сможет применить эту модель для выявления на ранней стадии пациентов из группы риска и помочь им вести здоровый образ жизни

* Из-за усреднения «вклада» деревьев может быть чрезвычайно сложно проследить логику в прогнозах.

путем консультаций и профилактических обследований. Из практики уже очевидны два признака, имеющие отношение к успешной диагностике: количество физических упражнений в неделю и вес. Теперь аналитик данных должен создать надежную модель, которая поможет достоверно прогнозировать, кто из пациентов попадет в группу риска.

Чего ожидать от k-NN?

k-NN анализирует вероятность. Метод заключается в вычислении расстояния между новой точкой данных и уже существующими. И поскольку существующие точки данных представляют собой ранее диагностированных пациентов, мы можем сгруппировать их в две категории: 1) страдающие диабетом и 2) здоровые. Затем новая точка данных (в нашем случае — новый пациент) будет классифицирована в соответствии с окружающими пациентами. Именно здесь мы наблюдаем основное допущение этого алгоритма: k-NN допускает, что даже *неизвестные* особенности пациентов будут схожи при условии, что схожи некоторые известные особенности.

Построение алгоритма k-NN

- 1. Выберите для вашего алгоритма число k — количество ближайших соседей.** Важно сначала установить, сколько соседних точек данных в нашем тренировочном наборе мы хотим проанализировать, чтобы новая точка данных была успешно классифицирована. k-NN анализирует расстояние между нашей новой точкой данных и существующими точками вокруг нее и классифицирует новую точку данных в соответствии с категорией (здесь либо (1) страдающий диабетом, либо (2) здоровый), которая представлена наибольшим числом соседей. Например, если мы хотим классифицировать наши новые точки данных посредством анализа пяти ближайших точек данных, то мы определим значение k как 5*.
- 2. Измерьте (евклидово) расстояние между новой точкой данных и всеми существующими точками.** Раз мы сказали, что

* Обычно используемое значение k в k-NN равно 5 и является числом по умолчанию для многих инструментов анализа данных.

k равно 5, то нам нужно определить пять соседей, ближайших к нашей точке данных. Для этого мы должны сначала измерить расстояние от нашей новой точки данных до всех точек, которые у нас уже есть.

В науке о данных расстояние может быть измерено несколькими способами. Обычно используется наиболее естественное расстояние — евклидово, то, что многие из нас изучали в школе. Евклидово расстояние — это длина отрезка прямой между двумя точками. Она измеряется путем нахождения разности в координатах двух точек для каждой оси (например, $X_2 - X_1$), затем их возведения в квадрат, суммирования результирующих значений и наконец извлечения квадратного корня.

Например, если P_1 — наша первая точка данных, а P_2 — вторая, как показано на графике (рис. 6.4), то евклидово расстояние будет измеряться по формуле:

$$\text{евклидово расстояние} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

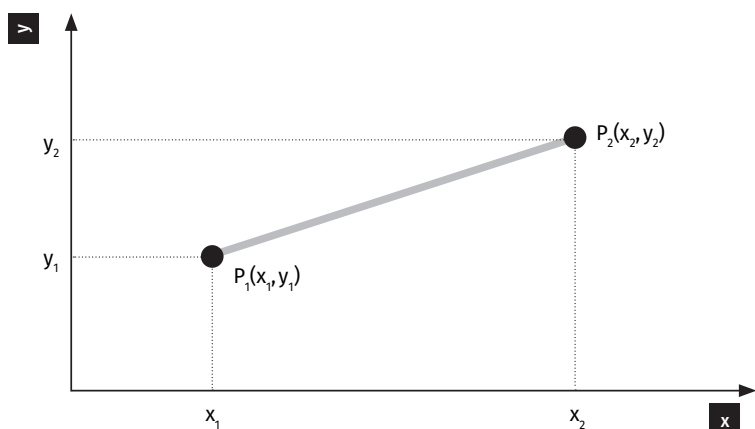


Рис. 6.4. Измерение евклидова расстояния

Вы, наверное, помните теорему Пифагора из школьной программы, и это точно такой же принцип. Две точки на этом графике являются двумя вершинами прямоугольного треугольника. Гипотенузу можно найти, сначала определив квадрат ее длины: для этого сложите квадраты двух других сторон.

3. Подсчитайте количество точек данных в каждой категории.

После того как вы нашли евклидово расстояние между новой точкой данных и каждой из старых точек данных, вы должны ранжировать эти расстояния в порядке возрастания. На данном этапе будет легко определить k -ближайших соседей — это просто первые пять пунктов в вашем списке. Визуально мы можем обвести ближайшие точки данных следующим образом (рис. 6.5):



Рис. 6.5. Алгоритм k -NN

4. Отнесите точку данных к категории с наибольшим количеством соседей. Мы видим, что для нашей новой точки данных есть три ближайших соседа в категории 1 — страдающих диабетом и только два ближайших соседа в категории 2 — здоровых. Поскольку в категории 1 больше близких соседей, мы отнесем новую точку данных к этой же категории, что означает, что этот конкретный пациент, учитывая его вес и количество выполняемых физических упражнений, подвержен риску развития диабета второго типа. Так мы классифицировали новую точку, модель готова*.

* Это показывает, почему общепринятой практикой является выбор $k = 5$: тестирование нечетного числа соседей помогает избежать ситуаций, когда категории становятся связанными с равным числом «ближайших соседей» (например, 2 и 2 ближайших соседа, если $k = 4$).

Многомерные пространства

Что происходит, если у нас более двух переменных для описания точек данных? Что, если в дополнение к весу и количеству физических упражнений у нас также была бы информация о возрасте пациентов и их среднем ежедневном потреблении калорий? Поскольку у нас есть несколько переменных, мы уже не можем рисовать двумерную диаграмму рассеяния. Вместо этого нам нужна четырехмерная диаграмма.

Представить или визуализировать точечную диаграмму 4D практически невозможно, но хорошая новость в том, что k -NN будет работать независимо от этого, поскольку алгоритм основан на подобии, зависящем от расстояния, — и формула, которую мы ввели для расстояния, может быть переписана для любого количества измерений. Просто будет больше элементов под квадратным корнем.

Тестирование

Как бы ни был хорош k -NN для создания точных прогнозов, важно отметить, что результаты применения этого метода не всегда будут правильными. Это совершенно нормально — всегда будет несколько неверных прогнозов и ни один алгоритм не сможет всегда давать правильные ответы. Ключ к созданию хорошей модели состоит в том, чтобы проверить ее несколько раз, изменяя функции (в нашем случае — значение k), пока вы не найдете лучшее решение для своей задачи.

Плюсы и минусы использования алгоритма k -NN

Алгоритм k -NN часто является правильным выбором, потому что он интуитивно понятен и, в отличие от наивного байесовского классификатора, как мы увидим ниже, не разрешает допущения о данных. Однако основным недостатком k -NN является то, что вычисление занимает очень много времени. Необходимость вычислять расстояние до каждой точки в массиве данных чревата тем, что чем больше у вас точек, тем медленнее k -NN будет работать.

Наивный байесовский классификатор

Наивный байесовский классификатор назван в честь теоремы Байеса, которая позволяет математикам выражать вероятность событий таким образом, что любые вновь открытые доказательства могут быть легко включены в алгоритм для динамического обновления значения вероятности. Это увлекательный алгоритм, потому что он позволяет видеть сквозь созданную нашим разумом иллюзию и проливает свет на реальное положение дел.

Чтобы лучше понять наивный байесовский классификатор, мы должны сначала взглянуть на теорему Байеса и ее уравнение. Как только мы разберем эти понятия, переход от теоремы к алгоритму классификации пройдет гладко.

Полицейские проверки и теорема Байеса

Вас когда-нибудь останавливал сотрудник полиции для проверки на алкоголь? Это распространено в Австралии в пятницу и субботу вечером, когда люди возвращаются домой с вечеринок, — австралийская полиция известна привычкой оцепить главную дорогу в самом оживленном месте. Любой, кто едет по этой дороге, независимо от манеры вождения должен остановиться для проверки уровня алкоголя в крови. Это быстрый процесс, так как вам даже не нужно выходить из автомобиля, и он помогает сотрудникам полиции убрать пьяных водителей с улиц. Мы собираемся использовать этот пример, чтобы лучше понять теорему Байеса.

Давайте поговорим об алкотестере. Предположим, что это устройство разработано очень хорошо и безошибочно выявляет всех пьяных водителей. В конце концов, в том его основное предназначение. Но алкотестер неидеален и будет регистрировать ложное пьянство в 5% случаев. Это означает, что из 100 *трезвых* человек он будет ошибочно считать пьяными пятерых (такие результаты называются ложно-положительными). То есть тестирование алкотестером даст положительный результат, хотя на самом деле эти люди не будут находиться в состоянии опьянения.

А теперь представьте, что полицейский только что проверил алкотестером случайного водителя и прибор показывает, что водитель пьян. Какова вероятность того, что он или она действительно выпили?

Импульсивный ответ был бы 95%. Но правильный ответ на самом деле — около 2%. Как так? Здесь пригодится теорема Байеса.

Предположим, что на каждую 1000 водителей на дороге приходится только один человек, который ведет машину в нетрезвом состоянии. Если полицейские протестируют 1000 водителей, они получат следующие результаты:

- 1 водитель, который действительно пьян, будет обнаружен непременно;
- из оставшихся 999 водителей 5% будут сочтены пьяными, то есть $5\% \times 999 = 49,95$ водителя (не беспокойтесь о десятичной запятой в числе водителей — мы всегда можем округлить этот пример до 100 000 водителей, чтобы результат был целым числом).

В этом примере алкотестер выявил в общей сложности $1 + 49,95 = 50,95$ пьяного водителя. Таким образом, вероятность того, что любой из этих водителей действительно пьян, $1/50,95 = 0,0196\% \approx 2\%$. Мы можем проиллюстрировать это в таблице 6.1:

Таблица 6.1

	Действительно пьяные (Actually drunk)	Действительно непьяные (Actually not drunk)	Итого (Total)	
Протестировано	1	999	1000	
Выявлено	$\times 100\%$	$\times 5\%$		
	= 1	= 49,95	50,95	$P = 1/50,95$ = $1,96/100$ = 1,96%

Удивлены? Вы не одиноки. Теорема Байеса до сих пор озадачивает меня всякий раз, когда я сталкиваюсь с примером ее применения к реальной жизненной ситуации. Поразительно, как часто мы делаем

поспешные выводы о том, что нам показывают, вместо того чтобы рассмотреть общую картину*.

Формула Байеса

Теперь давайте посмотрим на формулу Байеса. Вот обозначения, которые будут использоваться в этом примере:

P (пьяный);
 P (пьяный | положительно);
 P (положительно | пьяный);
 P (положительно),

где P обозначает вероятность, а вертикальная черта — условную вероятность.

Каждый из перечисленных элементов имеет математическое название. P (пьяный) — вероятность того, что случайно выбранный водитель будет пьян. В байесовской статистике эта вероятность называется *априорной вероятностью*. Если мы вспомним наши первоначальные предположения, то можем вычислить априорную вероятность как P (пьяный) = $1/1000 = 0,001$.

P (пьяный | положительно) — условная вероятность того, что при положительном результате алкотестера (когда устройство определило, что человек за рулем находится в состоянии алкогольного опьянения) водитель действительно окажется нетрезв. Эта вероятность называется *апостериорной вероятностью*, она нас интересует в расчете.

P (положительно | пьяный) — условная вероятность того, что, когда водитель фактически пьян, алкотестер отреагирует положительно. Ее называют *функцией правдоподобия*. В нашем случае любой по-настоящему пьяный водитель всегда распознается прибором, а значит, P (положительно | пьяный) = 1.

* См. видео Джулии Галеф «Визуальное руководство по байесовскому мышлению» для некоторых неожиданных реальных приложений теоремы Байеса. www.youtube.com/watch?v=BrK7X_XlGB8.

P (положительно) — вероятность того, что у любого случайно выбранного водителя окажется положительный результат на алкотестере. Это *предельное правдоподобие*, и в нашем примере оно рассчитывается как $P(\text{положительно}) = 50,95/1000 = 0,05095$.

Не волнуйтесь, вы не должны помнить все эти названия, но в один прекрасный день можете встретить их — и тогда припомните наш пример с алкотестером. А теперь, когда все приготовления завершены, можно ввести формулу Байеса:

$$P(\text{пьяный} \mid \text{положительно}) = \frac{P(\text{положительно} \mid \text{пьяный}) \times P(\text{пьяный})}{P(\text{положительно})}$$

Подставив числа, получим следующее:

$$1 \times 0,001/0,05095 = 0,0196 = 1,96\%.$$

Хотя это уравнение может показаться сложным, на самом деле оно понятно на уровне интуиции. Если вы не уверены, просто повторите шаги, с помощью которых мы рассчитали этот тип вероятности, для начала используя табличный метод, и вы увидите, что мы выполнили точно такие же вычисления, как предложено формулой Байеса. Разница лишь в том, что наши исходные данные были приведены к 1000 водителей (вместо 0,001 у нас был 1, а вместо 0,05095 — 50,95).

Дополнительные свидетельства

Мы так увлеклись теоремой Байеса, что совсем забыли о наших полицейских. Они проверяют человека алкотестером, прибор считывает данные как положительные, но вероятность того, что человек действительно пьян, всего лишь 2%. Что делать копам?

На этом этапе они могут прибегнуть к более точным методам проверки (например, взять анализ крови на уровень алкоголя) или выбрать гораздо более простое решение: дополнительное тестирование содержимого выдоха. Давайте посмотрим, чем это может быть полезно.

Мы знаем, что из 1000 протестированных водителей 50,95 были признаны пьяными. Мы также знаем (для целей этого примера), что только один из них *на самом деле* пьян. Тестируя каждого из 50,95 «пьяного» водителя вторично, можно применить ту же логику, что и раньше, — алкотестеры обнаружат:

- одного водителя, который на самом деле пьян;
- 5% остальных 49,95 водителя в нетрезвом состоянии, то есть $5\% \times 49,95 = 2,4975$ водителя.

Таким образом, $1 + 2,4975 = 3,4975$ водителя будут сочтены пьяными во второй раз.

Как видно, мы сужаем область поиска результатов, и теперь вероятность того, что водитель, чей результат во *втором* испытании алкотестером был положительным, действительно пьян, равна: $1/3,4975 = 28,59\%$.

Результат по-прежнему кажется низким? Тогда почему бы еще раз не протестировать остальных водителей? Применяя ту же логику, мы получим следующие результаты проверки в третьем раунде:

- как всегда, один действительно пьяный водитель будет обязательно обнаружен;
- из оставшихся 2,4975 водителя 5% будут сочтены пьяными, то есть $5\% \times 2,4975 = 0,124875$ водителя.

Теперь только $1 + 0,124875 = 1,124875$ водителя сочтен пьяным. Таким образом, вероятность того, что водитель с положительным результатом третьего тестирования алкотестером действительно был пьян, равна: $1/1,124875 = 88,89\%$.

Вот так намного лучше. На этом этапе сотрудники полиции могут приказать водителям, тестирование которых дало положительный результат, выйти из автомобилей. Четвертый тест будет еще более точным, и вероятность совпадения его результатов и состояния водителя возрастет до более 99%. Вы вполне можете выполнить этот расчет в свободное время. Чтобы не потеряться в числах, используйте таблицу 6.2 в качестве руководства:

Таблица 6.2

	Действительно пьяные	Действительно непьяные	Итого	
Протестировано	1	999	1000	
Выявлено	$\times 100\%$ $= 1$	$\times 5\%$ $= 49,95$	50,95	$P = 1/50,95$ $= 1,96/100$ $= 1,96\%$
Тест 2	$\times 100\%$	$\times 5\%$		
Выявлено	$= 1$	$= 2,4975$	3,4975	$P = 1/3,4975$ $= 28,59\%$
Тест 3	$\times 100\%$	$\times 5\%$		
Выявлено	$= 1$	$= 0,12487$	1,12487	$P = 1/1,12487$ $= 88,89\%$

Если вы хотите следовать формуле Байеса (в отличие от табличного метода), все, что вам нужно сделать, — это обновлять на каждом шаге исходные данные (априорная вероятность и предельное правдоподобие), включая результаты предыдущего шага. Это математический способ сказать, что вы получили новые доказательства и хотите использовать их для уточнения вашего существующего (априорного) взгляда на мир.

Пример с использованием алкотестера мне очень нравится, поскольку он иллюстрирует две вещи:

1. Нужно с самого начала учитывать априорные знания (в данном примере, что только один из 1000 водителей фактически пьян). Игнорирование общей картины может привести к поспешным и зачастую неправильным выводам.
2. Исходные данные в формуле Байеса нужно обновлять по мере поступления новых. Только так можно добиться того, чтобы общая картина всегда оставалась актуальной. Иногда нам, возможно, придется активно искать новые сведения, чтобы получить более точные результаты.

Так что читатели больше не должны удивляться, что сотрудники полиции иногда просят водителей дышать в контрольные устройства более одного раза.

Это был краткий экскурс в мир байесовской статистики. Вооружившись теоремой Байеса, мы теперь готовы перейти к наивному байесовскому классификатору.

Почему он наивный?

Наивный байесовский классификатор основан на сильном, *наивном* допущении независимости признаков: все характеристики массива данных не зависят друг от друга. На самом деле было бы наивным так полагать, поскольку для многих массивов данных может быть выявлен уровень корреляции содержащихся в них независимых переменных. Несмотря на это наивное предположение, наивный алгоритм Байеса хорошо зарекомендовал себя во многих сложных приложениях, таких как программа для обнаружения спама в электронной почте.

Использование наивного байесовского классификатора

К нам обратился винодел из Калифорнии. Погода на Западном побережье тогда установилась капризная, и винодел опасался за качество будущего вина. Ему нужна была помощь в прогнозировании шансов его продукции возглавить региональный список лучших вин урожая того года.

Можно сказать, для нашего винодела многое было поставлено на карту. Хорошая новость состояла в том, что у него имелись некоторые данные для нас!

Винодел обнаружил, что на протяжении многих лет две независимые переменные — *продолжительность солнечного сияния* и *количество осадков* — оказывают положительное влияние на виноградные лозы и, соответственно, на вкус его вин, а значит, повышают шансы на успех. С тех пор ему удалось усовершенствовать процесс выращивания винограда и тем самым улучшить качество своей продукции.

Основываясь на своих предыдущих победах и поражениях, винодел разделил имеющиеся у него данные на две категории: «победитель» и «проигравший». Мы можем визуальнo представить их так (рис. 6.6):

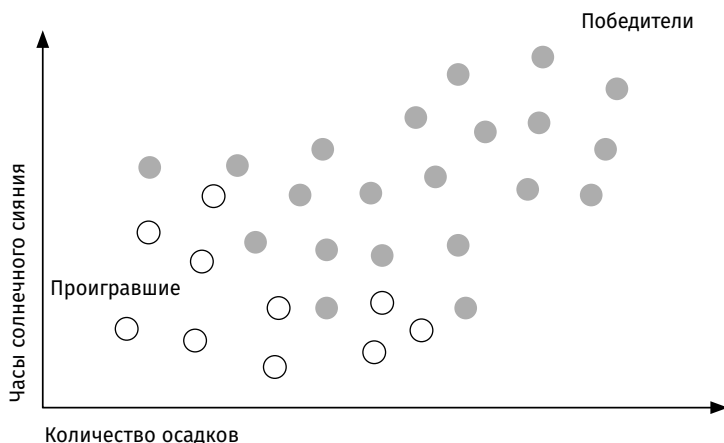


Рис. 6.6. Качество урожая винограда в зависимости от продолжительности солнечного сияния и количества осадков (1)

Здесь значение по оси x — миллиметры осадков, а значение по оси y — часы солнечного сияния. Белая категория — «проигравший», а серая — «победитель». Теперь мы можем помочь виноделу проанализировать шансы на успех вина из урожая этого года, основываясь на количестве осадков и продолжительности солнечного сияния. Предположим, что в период созревания конкретного урожая выпало 601,98 мм осадков и что на это время пришлось 3543 часа солнечного сияния. Используя эту информацию, мы можем построить график рассеяния для урожая этого года, и наивный классификатор Байеса поможет нам определить, в какую категорию попадет урожай этого года (рис. 6.7).

Построение наивного байесовского классификатора

Наивный байесовский классификатор использует переменные нашей точки данных, чтобы отнести ее к наиболее подходящему классу. Вот как это работает.

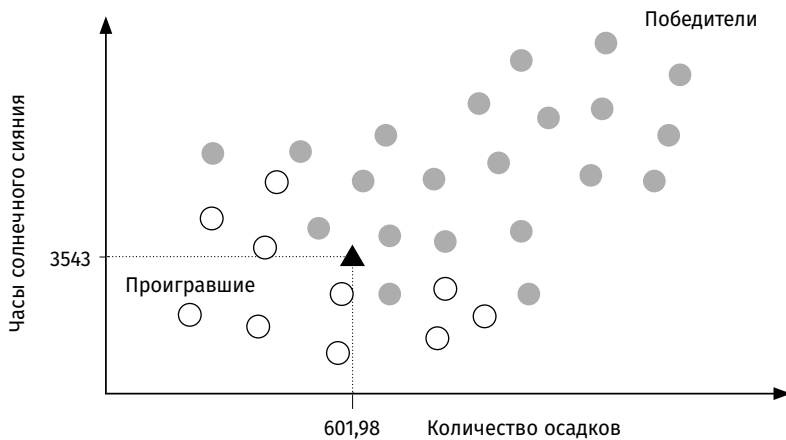


Рис. 6.7. Качество урожая винограда в зависимости от продолжительности солнечного сияния и количества осадков (2)

Шаг 1: установите априорную вероятность. Здесь мы хотим узнать вероятность того, что отдельная точка данных принадлежит к категории из нашего тренировочного набора. Учитывая размер выборки и количество проигравших и выигравших вин, какова вероятность того, что новое вино попадет в категорию победителей?

В этот момент нам нужно предположить, что мы ничего не знаем об урожае — нам неизвестно, сколько времени виноград провел на солнце и сколько выпало дождей. Так что лучшее, что мы можем сделать, — это взять количество победителей из наших предыдущих (априорных) данных (отсюда и название: *априорная* вероятность) и разделить его на общее число точек данных:

$$P(\text{победитель}) = \frac{\text{Количество победителей}}{\text{Общее количество наблюдений}} = 20/30, \text{ или } 0,667.$$

Шаг 2: вычислите предельное правдоподобие. Предельное правдоподобие относится к вероятности того, что новая точка данных находится в непосредственной близости от области, куда фактически попадает рассматриваемый вариант. Обычно или необычно для урожаев получать такое же количество солнечного света и осадков, как получил наш урожай? Это *условие подобия*

представляет собой область вокруг нашей точки данных, которая будет выглядеть примерно так на диаграмме рассеяния* (рис. 6.8).

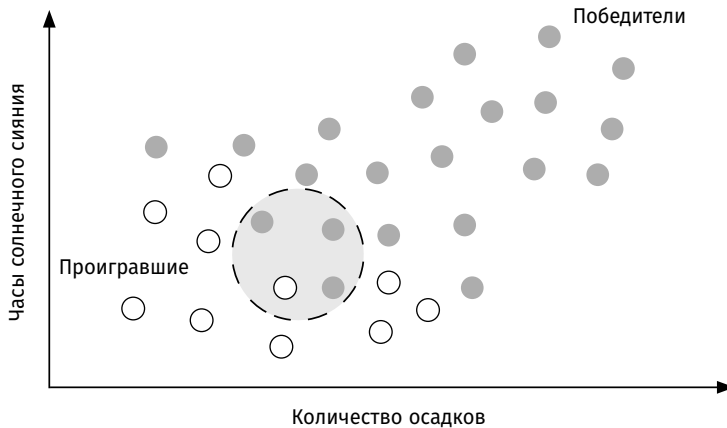


Рис. 6.8. Расчет предельного правдоподобия для наивного байесовского классификатора

Радиус круга мы выбираем произвольно; это параметр, который мы можем настраивать, чтобы влиять на эффективность алгоритма.

Таким образом, точки данных, содержащиеся в пределах нашей окружности, считаются одинаковыми. Эти вина сделаны из винограда, получившего примерно такое же количество солнечного света и воды, что и наш сегодняшний урожай. Допустим, что наш круг включает четыре точки данных. Чтобы найти вероятность того, что новая точка данных (X) попадет в круг, нам нужна следующая формула:

$$P(X) = \frac{\text{Аналогичные наблюдения}}{\text{Общее число наблюдений}} = \frac{4}{30}, \text{ или } 0,133.$$

Обратите внимание, что это значение не изменится в течение всего времени нашего анализа, поэтому его достаточно рассчитать только один раз.

* Имейте в виду, что мы на самом деле не размещаем наш текущий урожай на диаграмме рассеяния — это позволяет избежать путаницы, когда мы начинаем подсчет точек данных в круге. Вместо этого мысленно представьте, что он находится в середине круга.

Шаг 3: вычислите функцию правдоподобия (рис. 6.9). Как мы помним из теоремы Байеса, *функция правдоподобия* является условной. Какова вероятность того, что точка данных в нашем массиве данных попадет в круг, который мы определили, *учитывая*, что она уже принадлежит к категории победителей?

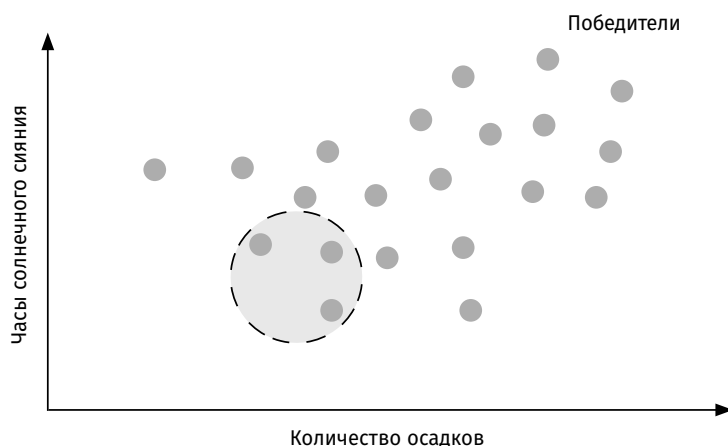


Рис. 6.9. Расчет функции правдоподобия для наивного байесовского классификатора

Чтобы найти функцию правдоподобия, нужно просто разделить количество аналогичных наблюдений в этой категории (в данном случае их три) на общее количество точек данных в категории:

$$\begin{aligned}
 P(X \mid \text{победитель}) &= \\
 &= \frac{\text{Количество аналогичных наблюдений среди вин-победителей}}{\text{Общее количество победителей}} = \\
 &= 3 / 20 = 0,15.
 \end{aligned}$$

(В случае, если вы находите обозначение $P(X \mid \text{победитель})$ запутанным, подумайте о букве «X» как о требовании, чтобы урожай, который мы принимаем в расчет, обладал характеристиками (количество солнечного света и осадков), очень похожими на те, о которых сообщил винодел. Поэтому запись $P(X \mid \text{победитель})$ равносильна вопросу «Какова вероятность того, что у этого вина будут такие же характеристики, какие винодел наблюдал у вина-победителя?».)

Шаг 4: рассчитайте апостериорную вероятность (ту, которая нас интересует!). Теперь нам нужно определить вероятность того, что из нового урожая может получиться вино-призер, притом что оно должно иметь за этот год характеристики (часы солнечного сияния и количество осадков), выявленные виноделом. Этот шаг выполняется с помощью формулы Байеса:

$$P(\text{Победитель} | X) = \frac{P(X | \text{Победитель}) \times P(\text{Победитель})}{P(X)}.$$

Подставив числа, получим следующее:

$$0,15 \times 0,667 / 0,133 = 0,75 = 75\%.$$

Это говорит о том, что из винограда *любого* урожая, который попадает в область, определенную нашим кругом, можно с вероятностью 75% произвести вино-победитель. Поэтому существует 75%-ная вероятность того, что и *наш* урожай позволит произвести наилучшее вино.

Шаг 5: выведите апостериорную вероятность противоположного сценария (из нашего урожая получится вино, неспособное победить). Теперь мы можем оценить вероятность того, что наше будущее вино попадет в проигравшую категорию, по аналогичной формуле*:

$$P(\text{Проигравший} | X) = \frac{P(X | \text{Проигравший}) \times P(\text{Проигравший})}{P(X)}.$$

Обратите внимание: следует повторить шаги 1 и 3 для проигравшей категории. Подстановка вычисленных вероятностей дает следующий результат:

$$((1/10) \times (10/30)) / (4/30) = 0,1 \times 0,333 / 0,133 = 0,25 = 25\%.$$

* Если в массиве данных есть только две категории, можно получить вероятность второго результата из первого, так как эти вероятности должны в конечном итоге составить 1 (или 100%). Тем не менее полезно использовать уравнение для всех категорий, чтобы дважды проверить, что они в сумме дают 1, — это хороший способ проверить ваши результаты.

Это говорит о том, что *любой* урожай, который попадает в область, определенную нашим кругом, с вероятностью 25% превратится в проигравшее вино. Поэтому существует 25%-ная вероятность того, что из *нашего* урожая получится вино, которое не попадет в число победителей.

Шаг 6: сравните две вероятности. Теперь мы знаем, что вероятность того, что новая точка данных (с характеристиками, отмеченными виноделом) принадлежит либо к категории победителя, либо проигравшего, составляет 75% и 25% соответственно — а это означает, что точка данных с этими характеристиками попадет в категорию победителей. Таким образом, хотя вероятность того, что вино из нового урожая проиграет, по-прежнему составляет 25%, вероятность его выигрыша выше, и поэтому данные, которые мы исследуем (урожай этого года), будут помещены в категорию победителей.

Наивный байесовский классификатор хорош для:

- нелинейных задач, в которых классы не могут быть разделены прямой линией на точечной диаграмме;
- массивов данных, содержащих сильно различающиеся данные (в отличие от других алгоритмов, результаты наивного классификатора Байеса нельзя исказить сильно различающимися данными).

Отрицательной стороной использования наивного классификатора Байеса является то, что сделанные с его помощью наивные предположения могут привести к погрешностям.

Вероятностные и детерминированные классификаторы

Это была долгая экскурсия в мир, где царствует наивный байесовский классификатор, не так ли? Я рад, что мы совершили ее, так как это важный алгоритм. Итог: есть 75%-ная вероятность того, что из урожая винограда в этом году получится вино-призер, и 25%-ная — того, что вино, произведенное из собранного в нынешнем году винограда, может проиграть соревнование.

Вы заметили, что этот вывод концептуально отличается от результата использования алгоритма k-NN?

Если бы мы применили k-NN к этому примеру, то получили бы однозначный ответ: вино из винограда нового урожая было бы объявлено либо победителем, либо проигравшим.

Черно-белое решение k-NN не дает иных возможностей, потому что, в отличие от наивного байесовского классификатора, принадлежит к семейству детерминированных алгоритмов классификации.

Детерминированные модели, такие как k-NN, относят полученные данные к одному конкретному классу, в то время как вероятностные модели, такие как наивный байесовский классификатор, предсказывают распределение вероятностей по всем классам. Затем это распределение можно использовать, чтобы отнести данные к классу.

Когда вы будете изучать следующий раздел об алгоритме логистической регрессии, спросите себя: *является ли этот алгоритм детерминированным или вероятностным классификатором?* Я скажу, правы ли вы, в конце раздела.

Логистическая регрессия

Несмотря на название, логистическая регрессия на самом деле не является алгоритмом регрессии; это тип метода классификации. Он использует наши данные, чтобы предсказать шансы на успех в таких сферах, как, скажем, продажа продукта определенной группе людей, определение ключевых демографических показателей для просмотра вашей электронной почты, или во многих других областях, не связанных с бизнесом, — например, в медицине, когда на основе возраста, пола и результатов анализа крови пациента пытаются предсказать, будет ли тот страдать ишемической болезнью сердца.

Но сначала мы должны вернуться назад. Для начала очень важно понять принципы линейной регрессии, в которую уходит корнями логистическая регрессия. Существует два типа линейной регрессии, о которых мы должны знать:

1. **Простая линейная регрессия** позволяет проанализировать связь между одной зависимой и одной независимой переменными. Это особенно полезно для анализа того, как одна переменная реагирует на другую, например когда мы рассматриваем изменение уровня преступности на фоне динамики ВВП страны.
2. **Множественная линейная регрессия** дает возможность проанализировать связь между одной зависимой и двумя или более независимыми переменными. Она лучше всего подходит для анализа более сложных массивов данных и может быть использована в целях изучения, например, того, каковы наилучшие предикторы (возраст, черты личности или социальная вовлеченность) уровней тревоги, испытываемой при смене жилья.

Как работает линейная регрессия

Ниже приведен пример линейной регрессионной модели на точечной диаграмме, которая показывает заработную плату респондентов и годы их стажа. Наша зависимая переменная — на оси y , а независимая переменная — на оси x (рис. 6.10).

При простой линейной регрессии, которую мы здесь наблюдаем, через наши данные проводится линия, и таким образом

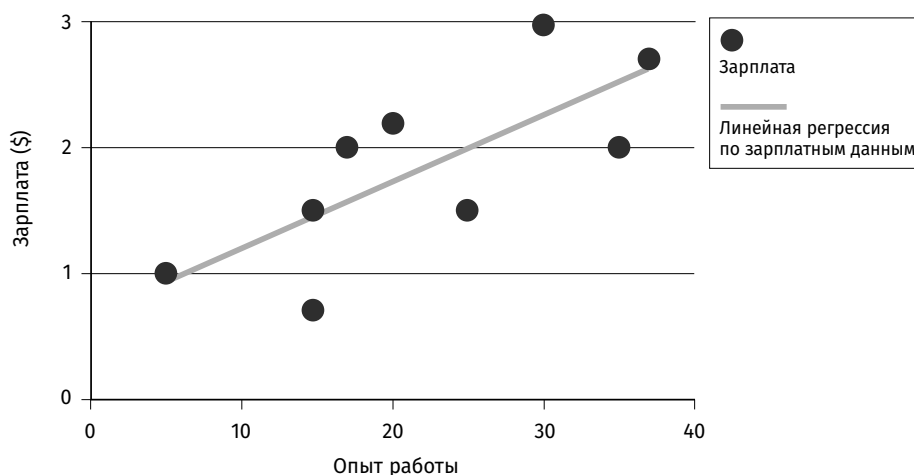


Рис. 6.10. Корреляционная диаграмма зависимости зарплаты респондентов от опыта

моделируются наши наблюдения. Это значит, что если мы будем знать опыт человека, то сможем спрогнозировать его зарплату. В то время как это хорошо работает для точечных диаграмм, где оси x и y содержат много значений, сложнее увидеть преимущества линейной регрессии для оси y только с двумя возможными значениями: 0 и 1. Это часто происходит, когда мы имеем дело с данными «да / нет», собранными из вопросов, на которые может быть дан один ответ из двух возможных. Вопросы типа «Вы купили этот продукт?», «Вы будете еще делать у нас покупки?» и «Есть ли у вас домашнее животное?» попадают в эту категорию, потому что требуют одного из двух ответов.

Ответы «да»/«нет»

Ответы «да»/«нет» являются *категориальными переменными*, то есть переменными с фиксированным числом ответов.

Работа с категориальными переменными. Можно ли найти регрессию для категориальных переменных? Да, можно. Давайте используем другой пример, чтобы проиллюстрировать это. Допустим, после e-mail-рассылки нашим клиентам мы хотим проанализировать уровень открываемости писем. На графике (рис. 6.11)

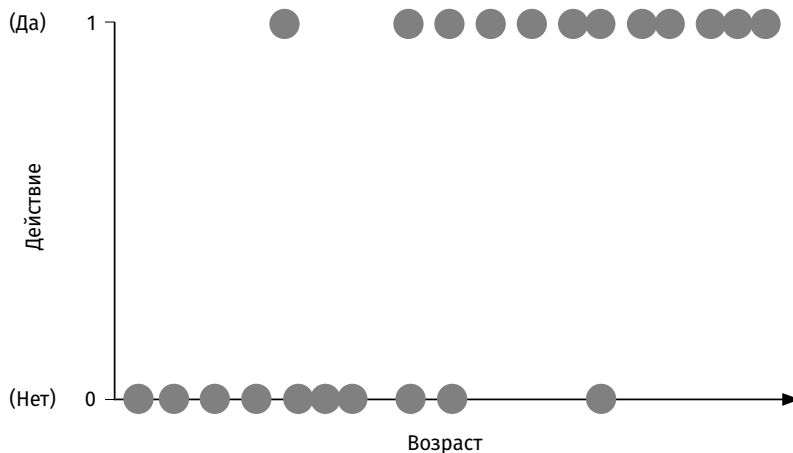


Рис. 6.11. Корреляционная диаграмма, показывающая, как респонденты открывают или не открывают наше сообщение в зависимости от их возраста

я показал, клиенты какого возраста открывают или не открывают наше электронное письмо. Значения «да»/«нет» были преобразованы в 1 и 0 соответственно.

На этом этапе мы можем задаться вопросом: что можно сделать со всем этим пространством между двумя значениями по оси y ? Как провести линию *регрессии* через график, который не показывает градиента изменений?

Но если мы посмотрим внимательнее, то увидим, что между значениями происходят постепенные изменения. На оси y значение 0 отклоняется влево по оси x , в то время как на оси y значение 1 больше отклоняется вправо по оси x . Это означает, что рассылка была хорошо принята пожилыми людьми. По мере увеличения значений на оси x (то есть с повышением возраста) рос стимул просмотреть наш e-mail. Это важный вывод, и теперь мы можем начать делать некоторые предположения о действиях, которые может предпринять человек определенного возраста.

Ось y нашего графика содержит значения 0 и 1. Мы также должны знать, что вероятности всегда имеют значения между 0 и 1. Таким образом, похоже, что линейная регрессия, которая проходит через интервал между этими значениями, даст нам информацию о *вероятности* того, откроет ли пользователь того или иного возраста наше электронное письмо (рис. 6.12).

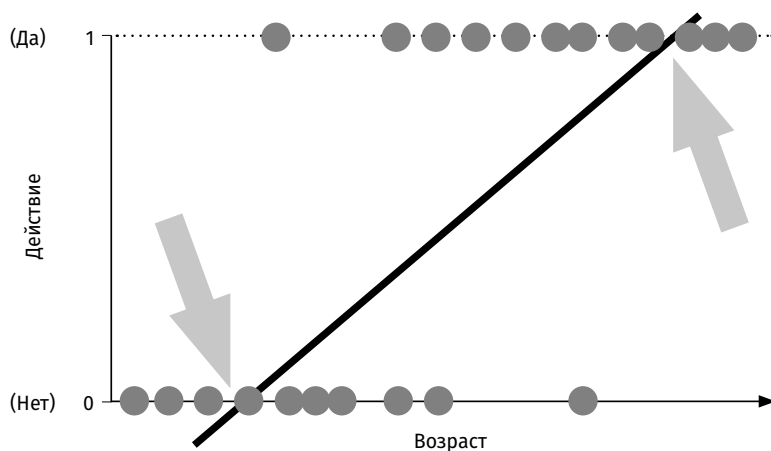


Рис. 6.12. Прямая линейная регрессия, показывающая вероятность открытия нашего сообщения респондентами в зависимости от их возраста

Сглаживание линии регрессии. Вы, возможно, заметили, что линия регрессии проходит по краям нашего графика. Это неидеально для вероятностей, так как они никогда не могут быть меньше 0 или больше 1, но могут быть только между двумя этими значениями.

Поэтому мы должны сократить части прямой, которые пересекают два значения 0 и 1. Как только линия линейной регрессии достигнет 0 или 1, она должна остаться на прямой и не продолжаться ниже или выше ее. Убедившись в этом, мы можем все так же использовать линию для создания предположений и быть уверенными, что наши результаты будут по-прежнему находиться в пределах вероятности. Первое, что нужно сделать, — обрезать несоответствующие части нашей линии (рис. 6.13):

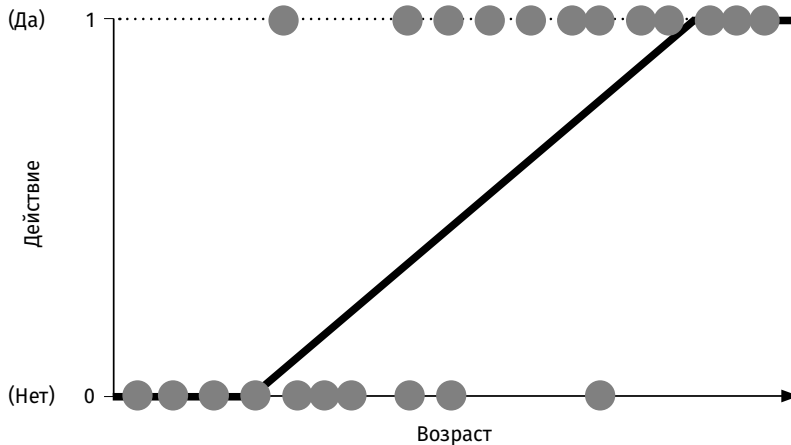


Рис. 6.13. Адаптированная линейная регрессия, показывающая вероятность открытия сообщения респондентами в зависимости от их возраста

Это хорошее начало, но есть более научный подход.

Математическая разработка логистической регрессии. График линейной регрессии может быть описан простым уравнением:

$$y = b_0 + b_1x.$$

Мы можем получить формулу логистической регрессии, если объединим приведенную выше формулу с так называемой сигмоидной

функцией* (функцией, график которой имеет форму S-образной кривой):

$$p = \frac{1}{1 + e^{-y}}.$$

После того как мы решим сигмоидную функцию для y и повторно вставим результат в первую формулу, мы получим:

$$\ln\left(\frac{p}{1-p}\right) = y = b_0 + b_1 \times x.$$

Эта формула преобразует наш график из прямой линии регрессии в функцию логистической регрессии (рис. 6.14):

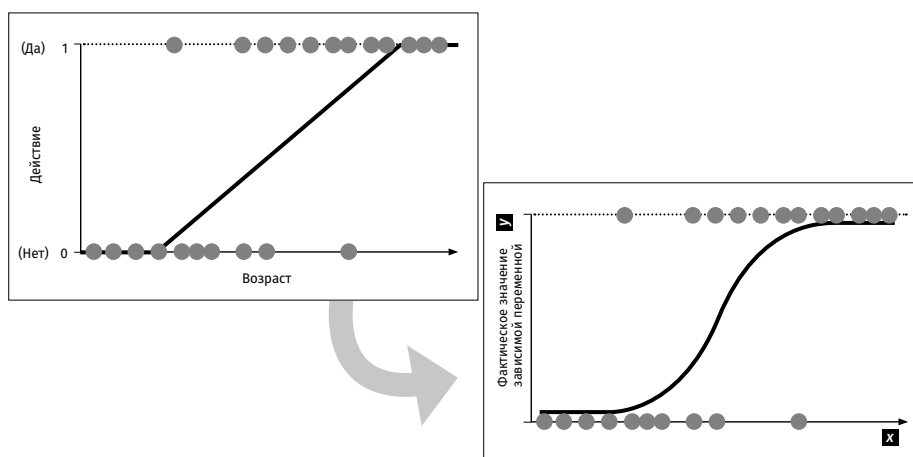


Рис. 6.14. Линейная регрессия, перенесенная на функцию логистической регрессии

Шаг 1: разберемся с элементами графика. Разобьем наш график на основные элементы (рис. 6.15):

* Не будем вдаваться в подробности и вместо этого станем опираться на интуицию. По этой причине некоторые уравнения будут даны без доказательств.

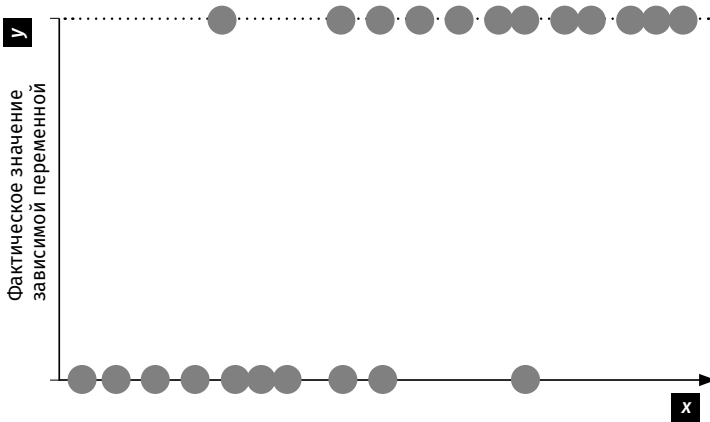


Рис. 6.15. График, содержащий категориальные переменные

Здесь ось x содержит независимую переменную, а ось y — зависимую с результатом «да»/«нет». Точки на графике — результаты, взятые из нашего массива данных.

Шаг 2: создание графика наклона для логистической регрессии. Мы сделаем это, подставив массив данных в формулу логистической регрессии и находя наиболее подходящие коэффициенты b_0 и b_1 :

$$\ln\left(\frac{p}{1-p}\right) = b_0 + b_1 \times x.$$

Это приводит к следующей кривой (рис. 6.16):

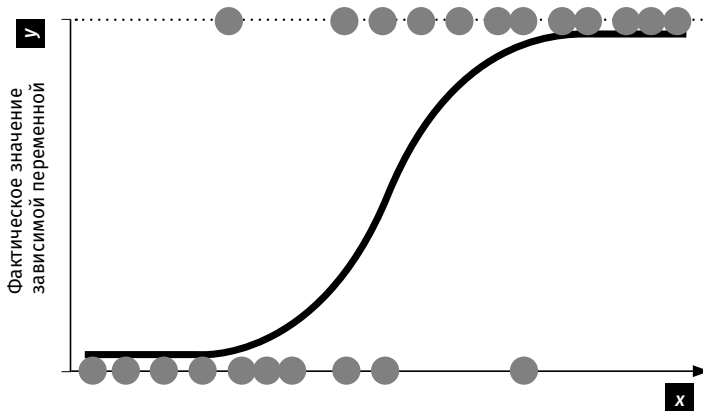


Рис. 6.16. Линия логистической регрессии

Эта кривая является наиболее подходящим графиком логистической регрессии для наших массивов данных. Как только мы проведем эту линию, можно стереть наблюдения из нашего графика, чтобы сосредоточиться на *самой линии* (рис. 6.17).

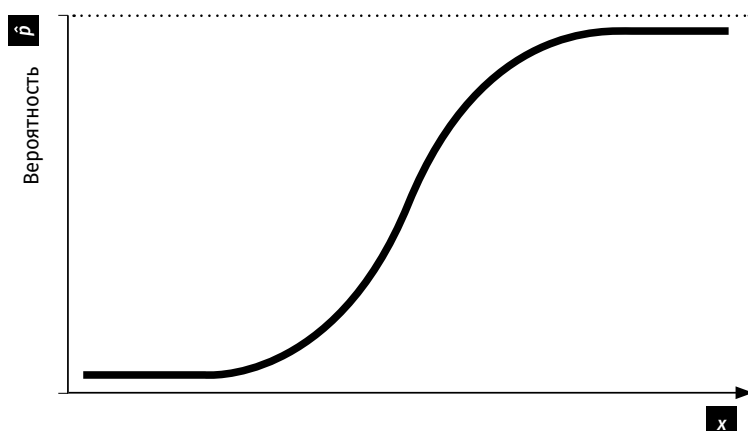


Рис. 6.17. Линия логистической регрессии (без наблюдений)

Обратите внимание, как изменилось обозначение оси y . Это потому, что мы можем использовать логистическую регрессию для прогнозирования вероятностей или *правдоподобия* того, что что-то произойдет. (На следующих страницах вы увидите символ $\hat{}$, например \hat{p} , — он означает предсказанные вероятности и называется крышечкой: \hat{p} — это p с крышечкой.)

Шаг 3: используйте график, чтобы сделать прогнозы для новых данных. Давайте вернемся к нашему примеру и предположим, что мы хотим определить вероятность открытия электронного письма людьми в возрасте 20, 30, 40 и 50 лет, учитывая, что у нас уже есть график логистической регрессии. Сначала мы спроецируем эти возрастные значения на кривую: проведем линии, параллельные оси \hat{p} , от каждой соответствующей точки на оси x до тех пор, пока они не достигнут линии регрессии. Это будут подходящие значения.

Затем мы проецируем эти значения влево, чтобы определить вероятность (рис. 6.18). Это означает проведение линии, параллельной

оси x , от установленного значения до тех пор, пока она не достигнет оси \hat{p} .

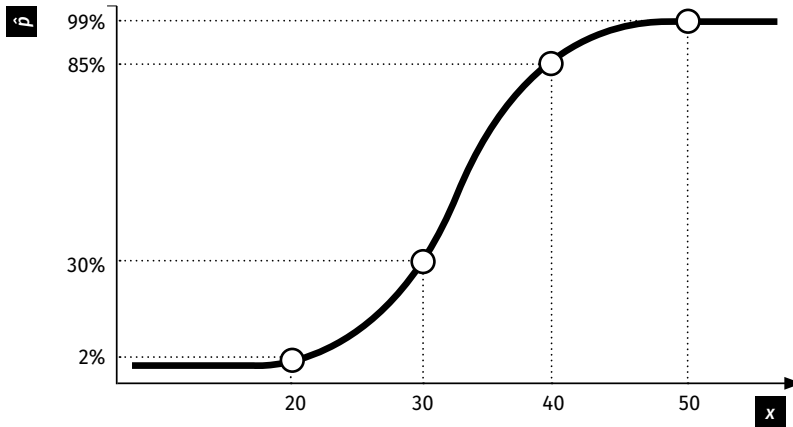


Рис. 6.18. Логистическая регрессия, включающая подходящие значения

Шаг 4: найти вероятность для каждого значения. Допустим (чисто гипотетически), что мы получили следующие результаты:

- возраст 20 лет дает вероятность 2% (или $\hat{p} = 2\%$);
- возраст 30 лет дает вероятность 30% (или $\hat{p} = 30\%$);
- возраст 40 лет дает вероятность 85% (или $\hat{p} = 85\%$);
- возраст 50 лет дает вероятность 99% (или $\hat{p} = 99\%$).

Шаг 5 (необязательно): установите ограничения. Итак, мы знаем, как получить вероятность \hat{p} для любой новой точки данных. Но как мы можем получить значение «да»/«нет»?

Хотя мы никогда не сможем быть абсолютно уверены в том, что произойдет, мы можем получить *предсказанное значение* для нашего фактического y (этот прогноз обычно обозначается \hat{y}) из нашей логистической регрессии.

Определить \hat{y} очень просто: выберите произвольный уровень на оси y между 0 и 1. Вы можете провести эту линию выше или ниже в зависимости от того, как много знаете о проблеме. Например, если вы продаете нишевый продукт, то, скорее всего, его купит меньше

людей, поэтому вы можете провести линию повыше, чтобы включить меньшее число потенциальных покупателей. Для этого примера давайте проведем линию прямо посередине, на 0,5 — это тоже самый распространенный подход (рис. 6.19).

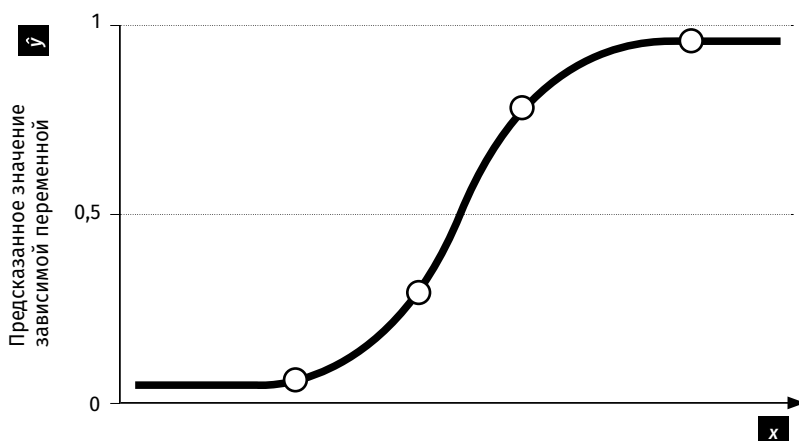


Рис. 6.19. Линия логистической регрессии (сегментированная)

Часть регрессии, которая находится ниже линии, установленной на уровне 0,5 (50%-ная вероятность), будет спроецирована на линию 0, чтобы дать $\hat{y} = 0$. Это означает, что если предсказанная вероятность открытия нашей электронной почты упадет ниже 50%, то мы можем предположить, что клиент, о котором идет речь, вероятно, не откроет наше электронное письмо. Все, что выше горизонтальной линии 0,5, будет проецироваться на линию 1, чтобы дать $\hat{y} = 1$.

Логистическая регрессия подходит для:

- 1) анализа вероятности заинтересованности клиента в вашем продукте;
- 2) оценки реакции клиентов на основе их демографических данных;
- 3) определения того, какая переменная является наиболее статистически значимой, то есть какая переменная оказывает самое большое влияние на зависимую переменную, значение которой мы хотим предсказать. Например, логистическая регрессия

может помочь нам определить, является ли статистика использования интернет-банкинга в течение последних шести месяцев более сильным предиктором оттока клиентов, чем сумма имеющихся у них сбережений.

Кластеризация

До сих пор мы говорили о классификации — о ситуации, когда мы всегда заранее знаем категории, в которые хотим сгруппировать или классифицировать новые точки данных. Теперь мы переходим к кластеризации, представляющей собой совершенно другое семейство алгоритмов.

Если вы не знаете, какими могут оказаться группы в результате анализа, следует использовать метод кластеризации. Методы кластеризации определенно сложнее, чем методы классификации, поскольку мы приступаем к решению задачи, не ведая, какие группы найдем.

Это не должно вас волновать. Именно поэтому мы сначала рассмотрели методы классификации. И то, к чему мы сейчас подходим, может быть гораздо более захватывающим процессом, ведь алгоритмы кластеризации позволяют нам использовать данные для того, чтобы обнаружить новые возможности и закономерности, чтобы выделить новые области, о которых мы, возможно, даже не подозревали, а не просто ответить на наш первоначальный вопрос.

В этом разделе мы рассмотрим алгоритм k -средних и алгоритм иерархической кластеризации. Эти два алгоритма во многом схожи, так как помогают разделить наши данные на статистически значимые группы.

Алгоритм k -средних

Алгоритм k -средних — одна из моих любимых моделей в науке о данных («модель» и «алгоритм» означают одно и то же). Хотя с ее помощью можно решать сложные задачи, метод k -средних легко понять; он основан на изящном интуитивном подходе.

Алгоритм k -средних обнаруживает статистически значимые категории или группы в нашем массиве данных. Это идеально подходит тогда, когда у нас есть две или более независимых переменных в массиве данных и мы хотим объединить точки данных в группы с похожими атрибутами. Например, k -средние могут помочь нам определить уровни подписки для массива данных членов киноклуба или показать комбинации групп интересов для интернет-магазина.

Построение алгоритма k -средних

Давайте проиллюстрируем алгоритм k -средних теоретическим примером.

Онлайн-ритейлер (наподобие Amazon или Alibaba) продает широкий спектр продуктов низкого и высокого класса. Некоторые пользователи платят за ежемесячную подписку, чтобы пользоваться услугами «премиум» от бесплатной доставки до раннего доступа к новым продуктам. Компания собрала большое количество данных (пол, возраст, годовой доход, история покупок, количество отдельных посещений

Почему бы не использовать здравый смысл?

Казалось бы, можно идентифицировать кластеры, просто разместив точки данных на точечной диаграмме, и *увидеть*, где окажутся наиболее значимые группы. Но использование алгоритма k -средних показывает нам группы, не видимые невооруженным глазом, и помогает провести линию (буквально) между группами, которые нам может быть трудно идентифицировать вручную.

Особенно это касается многомерных массивов данных. Для простоты мы свели задачу в нашем примере к двум переменным: годовой доход по оси x и оценка расходов по оси y . Но очень часто бывает так, что мы хотим выполнить кластеризацию по большему количеству измерений — трем, четырем, пяти, десяти, иногда даже больше. Поскольку визуализировать n -мерную диаграмму рассеяния невозможно, как мы видели в разделе о k -NN, алгоритм k -средних будет изящно решать n -мерную задачу, используя принципы, аналогичные тем, которые мы обсудим ниже.

веб-сайта в неделю и ежегодные расходы) о своих владельцах премиум-аккаунтов и получила «показатель расходов» для каждого из этих клиентов, вычисленный на основе комбинации числовых значений переменных.

Теперь компания хочет узнать, какого рода сегмент клиентов можно выделить среди существующих владельцев счетов. Благодаря такому сегментированию таргетинг предложений и маркетинга может стать значительно более эффективным, что позволит лучше удовлетворять индивидуальные запросы клиентов. В компании пока нет представления о том, какие группы они могут найти, — и это делает проблему подходящей для алгоритма кластеризации.

Чтобы разобраться с этим примером, мы приняли решение проанализировать годовой доход клиентов по отношению к их расходам.

1. Выберите число k кластеров. Вот наш массив данных (рис. 6.20):

Индекс	Номер клиента	Пол	Возраст	Годовой доход (тыс. \$)	Показатель расходов (1–100)
0	1	муж.	19	15	39
1	2	муж.	21	15	81
2	3	жен.	20	16	6
3	4	жен.	23	16	77
4	5	жен.	31	17	40
5	6	жен.	22	17	76
6	7	жен.	35	18	6
7	8	жен.	23	18	94
8	9	муж.	64	19	3
9	10	жен.	30	19	72
10	11	муж.	67	19	14
11	12	жен.	35	19	99

Рис. 6.20. Массив данных для клиентов компании электронной коммерции

Чтобы помочь понять роль интуиции при использовании алгоритма k -средних, предположим, что наш массив данных выглядит на точечном графике так, как показано на рис. 6.21 (в реальном мире данных было бы намного больше):

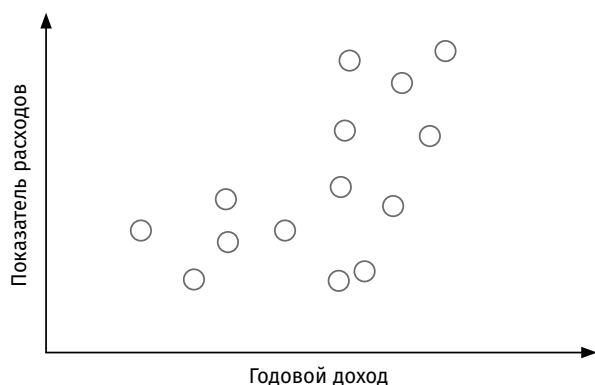


Рис. 6.21. Корреляционная диаграмма для двух независимых переменных

Наш первый шаг — выбрать количество кластеров (K), с которыми мы будем работать. Хотя одна из наших целей в k -средних — найти оптимальное количество кластеров, мы перейдем к этому позже (см. «Метод локтя и метрика ВКСК» ниже).

На данный момент нам нужно понять, как работает ядро алгоритма. Для этого возьмем случайное число кластеров и выберем два кластера в качестве отправной точки.

2. Выберите центроиды в случайных k точках. Так как мы ищем $k = 2$ кластера, мы случайным образом введем две точки данных (назовем их центроидами) в алгоритм*. Эти выбранные точки данных — место, откуда наш алгоритм начнет свое путешествие (рис. 6.22). Важно отметить, что такие точки данных не будут влиять на массив данных — это просто мнимые точки, которые алгоритм использует в процессе классификации.

3. Присвойте каждую точку данных ее ближайшему центроиду. Теперь центроид участвует в борьбе за территорию. Точки данных, ближайшие к центроиду А, будут относиться к нему на диаграмме, приведенной ниже, а точки данных, ближайшие к центроиду В, — относиться к В (рис. 6.23).

* Мы говорим «случайным образом» для простоты. Хотя при выборе начального местоположения центроидов необходимо помнить о некоторых подводных камнях, эта тема является более сложной и обычно учитывается алгоритмом.

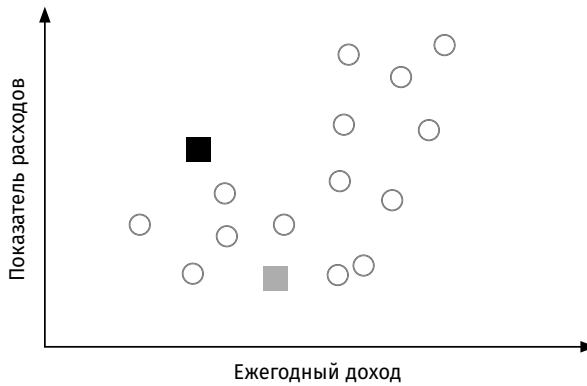


Рис. 6.22. k -средние, соответствующие случайным центроидам

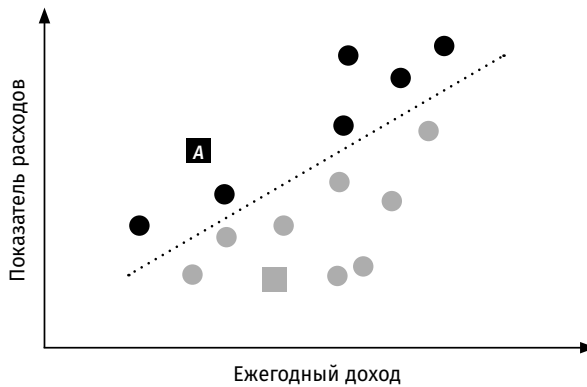


Рис. 6.23. Присвоение точек данных ближайшему центроиду

Пунктирная линия на диаграмме, равноудаленная от обоих центроидов, помогает увидеть, какие точки ближе к центроиду А, а какие — к центроиду В.

- 4. Определите и разместите новый центроид каждого нового кластера.** Теперь алгоритм вычислит «центр массы» для обоих сформированных кластеров. Затем он переместит каждый центроид в центр массы в соответствующем кластере (рис. 6.24).
- 5. Заново присвойте каждую точку данных новому ближайшему центроиду.** Борьба за власть продолжается — теперь, когда центроиды переместились в новые места, расстояния между центроидами и точками данных изменились. На этом этапе алгоритм уточнит, какие точки принадлежат А, а какие — В.

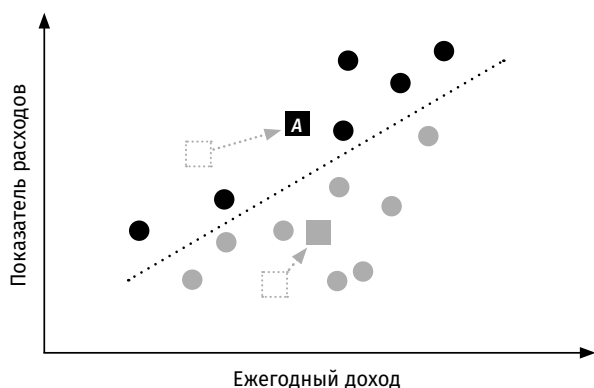


Рис. 6.24. Перерасчет местоположения центроидов

Прочерчивание равноудаленной линии между нашими центроидами А и В еще раз поможет нам визуализировать процесс (рис. 6.25).

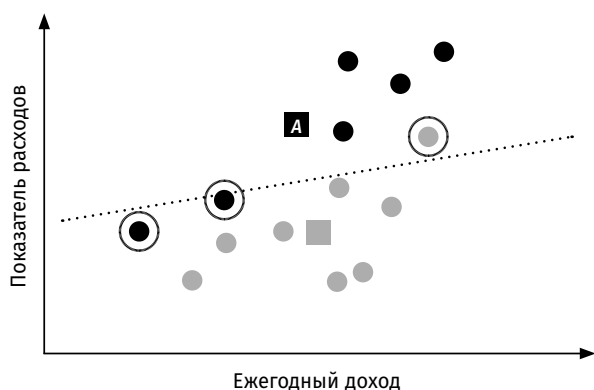


Рис. 6.25. Новое присвоение точек данных ближайшему центроиду

Как видно на диаграмме, три обведенные точки данных находятся на неправильной территории. Поэтому алгоритм заново присвоит эти точки ближайшим центроидам так, чтобы мы получили следующий результат.

Из рис. 6.26 видно, что центроиды вновь сместились — они не находятся в центре массы каждого кластера. Поэтому алгоритм повторит шаг 4, поместив центроиды в правильные места (рис. 6.27).



Рис. 6.26. Смещение центроидов

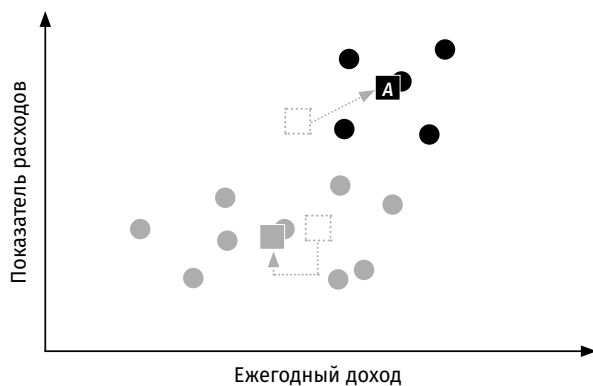


Рис. 6.27. Повторный перерасчет местоположения центроидов

Сразу после перемещения центроидов будет повторен шаг 5 и точки данных снова станут относиться к правильным центроидам. Алгоритм будет повторять шаги 4 и 5 до тех пор, пока не сойдется. Это произойдет тогда, когда на шаге 5 ни одной точке нельзя будет присвоить новый центроид и, следовательно, не будет смысла повторять шаги, поскольку ничего не изменится.

На рис. 6.28 показан наш конечный результат.

Как видно, центроиды проделали долгий путь от своих исходных точек. В итоге мы также получили две группы точек данных, которые, возможно, не были очевидны в начале нашего исследования (сравните с рис. 6.21).

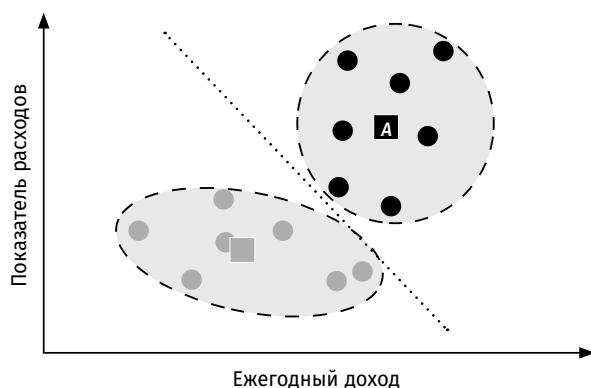


Рис. 6.28. k-средние, визуализация

Метод локтя и метрика ВКСК

Теперь, когда мы знаем, как работает алгоритм кластеризации k-средних, остается вопрос: как найти оптимальное количество кластеров (k) для использования? Этого можно добиться с помощью метода локтя.

Метод локтя помогает определить оптимальное количество кластеров. В качестве гипотетического примера рассмотрим другой массив данных, к которому мы применили k-средние, с $k = 3$ (рис. 6.29).

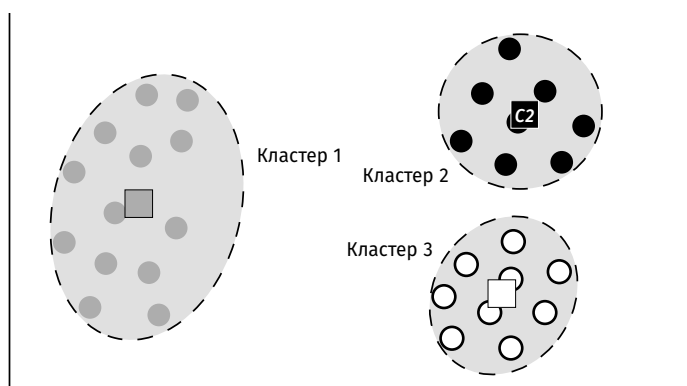


Рис. 6.29. Три кластера, идентифицированные с помощью k-средних

Хотя диаграмма выглядит весьма элегантно, она совсем не обязательно покажет нам оптимальные для решения этой проблемы группы. Разница между клиентами в кластерах 2 и 3, например, может быть недостаточно существенной для того, чтобы гарантировать их разделение на две группы. И наоборот, может оказаться более желательным разделить данные в кластере 1 на две или более группы.

Чтобы получить оптимальное число, мы должны оценить, как действует различное число кластеров, с помощью внутрикластерной суммы квадратов (ВКСК). Для трех кластеров, как указано выше, мы будем использовать следующую формулу ВКСК:

$$\begin{aligned} \text{ВКСК} = & \sum_{P_i \text{ в кластере 1}} \text{расстояние}(P_i, C_1)^2 + \\ & + \sum_{P_i \text{ в кластере 2}} \text{расстояние}(P_i, C_3)^2 + \\ & + \sum_{P_i \text{ в кластере 3}} \text{расстояние}(P_i, C_3)^2. \end{aligned}$$

Выглядит сложно? Не волнуйтесь, это на самом деле очень простое уравнение.

Три элемента

$$\sum_{P_i \text{ в кластере } X} \text{расстояние}(P_i, C_1)^2$$

просто рассчитываются для каждого кластера, и потом эти значения складываются для расчета ВКСК. Во втором элементе приведенного выше уравнения, например, мы рассматриваем только кластер 2. Сначала мы находим расстояние между каждой точкой и центром кластера (отмеченным его центроидом), а затем квадраты этих расстояний и наконец суммируем полученные значения.

Теперь (рис. 6.30) найдем ВКСК для алгоритма только с одним кластером

(k = 1):

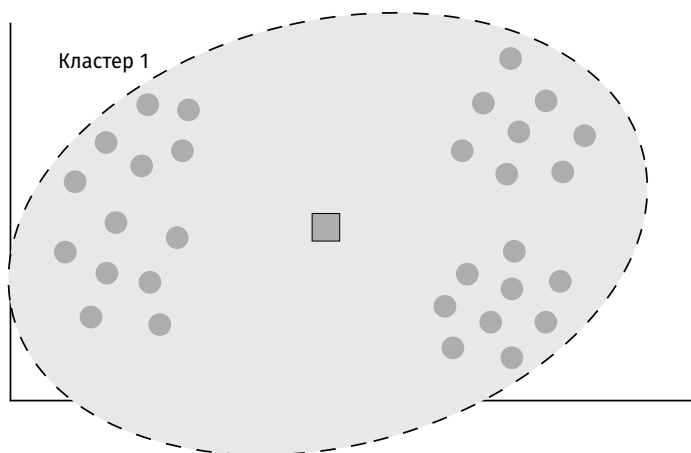


Рис. 6.30. Определение расстояния между точками данных в одном кластере

Здесь наши ВКСК потребовали бы, чтобы мы сложили квадраты всех расстояний между точками данных и центроидом. Мы получим большое значение ВКСК, потому что, как видно, центроид находится довольно далеко от некоторых точек.

Увеличим число кластеров до двух (рис. 6.31).

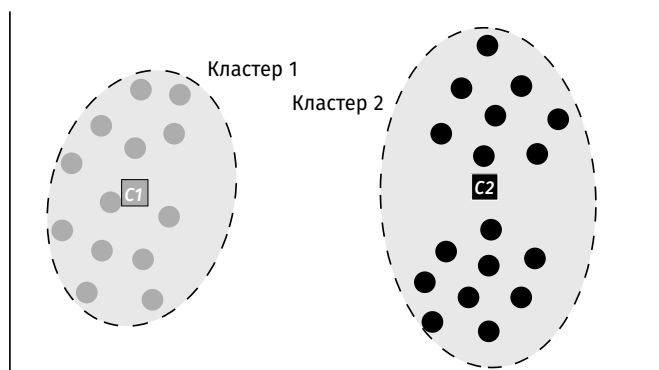


Рис. 6.31. Определение расстояния между точками данных в двух кластерах

Видно, что общее значение ВКСК будет меньше, чем вычисленное для одного кластера, поскольку расстояние между точками данных и центроидами соответствующих кластеров не так велико. Поэтому по мере увеличения числа кластеров значение ВКСК станет уменьшаться.

Каково же максимально возможное число таких кластеров?

У нас может быть столько кластеров, сколько точек в массиве данных. Например, 50 точек данных в массиве данных позволяют создать до 50 кластеров. Если количество кластеров совпадает с количеством точек данных, ВКСК равно 0, потому что центростид каждого кластера будет точно там, где точка данных.

Теперь, когда мы знаем, что значение ВКСК обратно пропорционально числу кластеров, мы можем изучить скорость изменения и получить наше оптимальное число. Нанесем эти значения ВКСК на график (рис. 6.32).

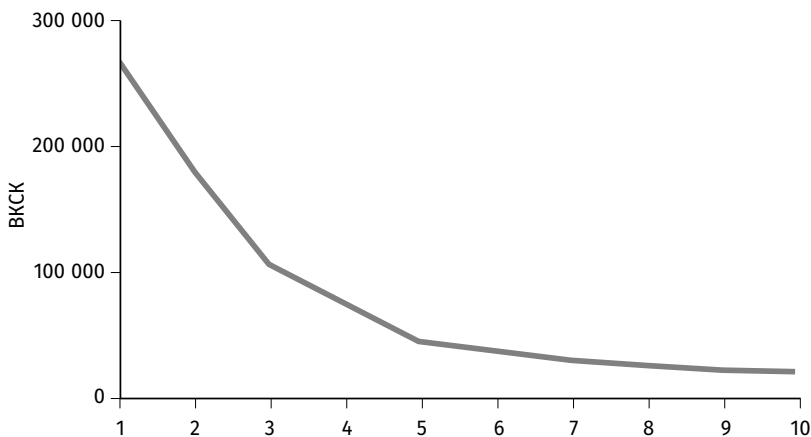


Рис. 6.32. Зависимость значений ВКСК от числа кластеров

Здесь мы видим, что после трех и пяти кластеров кривая ВКСК имеет излом. Чтобы определить оптимальное количество кластеров, все, что нам нужно сделать, — найти самый большой излом на графике; эта точка известна как «локоть». В случае выше было бы три кластера. Поэтому оптимальным количеством кластеров действительно будет $K = 3$.

Самое большое преимущество k -средних заключается в том, что их можно применять к массивам данных любого размера. Время исполнения растёт линейно с увеличением количества точек данных. Это не относится к другим алгоритмам, таким как иерархическая кластеризация (см. ниже), где время исполнения пропорционально квадрату числа точек данных. Разница не будет заметна на 10 или даже 100 записях, но представьте себе кластеризацию в 10 млн точек. Возможно, одним из самых больших недостатков k -средних является то, что даже при использовании метода локтя может быть трудно найти оптимальное количество кластеров.

Иерархическая кластеризация

Как и в случае с k -средними, мы будем использовать иерархическую кластеризацию, когда хотим сегментировать клиентов, но не знаем, сколько должно получиться групп или как наши данные могут быть разделены на части. В то же время, несмотря на то что оба алгоритма преследуют одну и ту же цель — идентифицировать данные по кластерам, иерархическая кластеризация и k -средние основаны на принципиально разных концепциях, и поэтому результирующие кластеры, вероятно, будут разными. Это еще одна причина узнать об обоих методах.

Существует два типа иерархической кластеризации (агломеративная и дивизивная), и они по существу являются двумя сторонами одной медали. Агломеративная иерархическая кластеризация использует подход «снизу вверх», работая с одной точкой данных и группируя ее с ближайшими точками данных поэтапно, пока все точки не будут собраны в один кластер.

Дивизивная иерархическая кластеризация работает противоположным образом. Она начинается сверху, где один кластер охватывает *все* наши точки данных, и прокладывает путь вниз, разделяя один кластер на части в зависимости от расстояния между точками данных. Процесс для обоих типов иерархической кластеризации записывается в так называемую дендрограмму.

Мы сосредоточимся здесь на агломеративной иерархической кластеризации, так как она наиболее часто используется.

Построение алгоритма агломеративной иерархической кластеризации

Шаг 1: сделайте каждую точку данных отдельным кластером. Прежде всего мы должны рассматривать наши отдельные единицы данных как кластеры.

Шаг 2: объедините два ближайших кластера. Возьмите два кластера, которые находятся ближе всего друг к другу, и объедините их. На рис. 6.33 одно это действие позволило сократить количество первоначальных шести кластеров до пяти. Сейчас мы повторим этот шаг, но с учетом этих пяти кластеров.

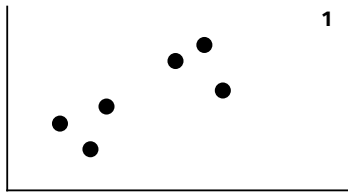
Повторяйте этот шаг, пока не останется только один кластер.

Определение расстояния

Даже если мы станем использовать евклидово расстояние (см. «Построение классификации случайного леса»), то, в отличие от ситуации с отдельными точками, расстояние между кластерами все еще будет неясно и должно быть точно определено. Вот несколько возможных вариантов измерения расстояния между двумя кластерами:

- A. Расстояние между их «центрами масс».
- B. Расстояние между двумя ближайшими точками.
- C. Расстояние между двумя самыми дальними точками.
- D. Среднее значение B и C.

Как правило, по умолчанию берется расстояние между центрами масс двух кластеров. Тем не менее ваш выбор здесь может значительно повлиять на конечные результаты — опирайтесь на свое внутреннее знание проблемы, чтобы сделать обоснованный выбор.



1

1. Рассматривайте каждую точку данных как одноточечный кластер

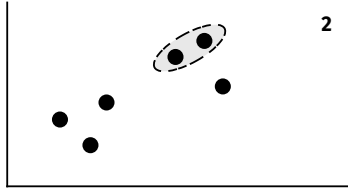
→ У вас есть 6 кластеров

2. Возьмите 2 ближайших элемента и сделайте их одним кластером

→ У вас есть 5 кластеров

3. Снова возьмите 2 ближайших элемента и сделайте их одним кластером

→ У вас есть 4 кластера



2

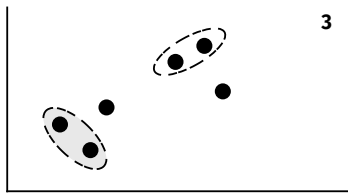
4. Повторите

→ У вас есть 3 кластера

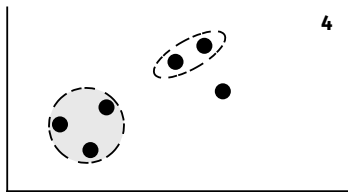
5. Повторите

→ У вас есть 2 кластера

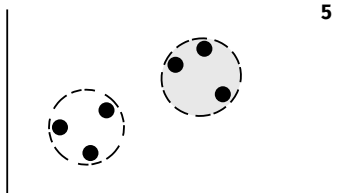
6. Повторяйте, пока не останется один кластер



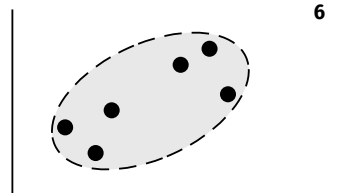
3



4



5



6

Рис. 6.33. Построение агломеративной иерархической кластеризации

Почему агломеративная кластеризация полезна?

Этот тип кластеризации содержит запись каждого этапа процесса: фиксируется порядок поглощения точек данных и расстояние между ними на древовидной диаграмме, известной как дендрограмма.

Что такое дендрограммы?

Дендрограмма расположит точки ваших данных (P1, P2, P3, P4) на оси x графика. Расстояния между точками данных представлены на оси y (рис. 6.34).

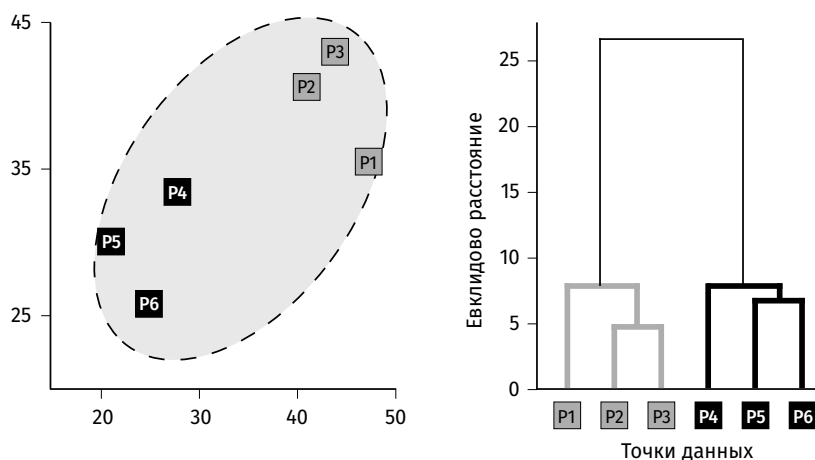


Рис. 6.34. Процесс агломеративной кластеризации, представленный на дендрограмме

Как видно, высота вертикальных линий зависит от расстояния между каждой точкой данных, а горизонтальные линии показывают порядок, в котором происходила кластеризация. Самые низкие горизонтальные линии представляют первые объединенные кластеры; дальнейший путь вверх показывает процесс группировки. В этом примере мы видим, что первые две кластерные точки — P2 и P3, затем P5 и P6. Затем P1 была кластеризована с P2 и P3, а P4 — с P5 и P6. Наконец, эти две группы (P1, P2, P3 и P4, P5, P6) были кластеризованы.

Шаг 3: установите порог. С помощью дендрограммы можно установить порог, который позволит узнать, какое число кластеров оптимально для нашего проекта. Нарисуем произвольный порог на нашей дендрограмме (рис. 6.35).

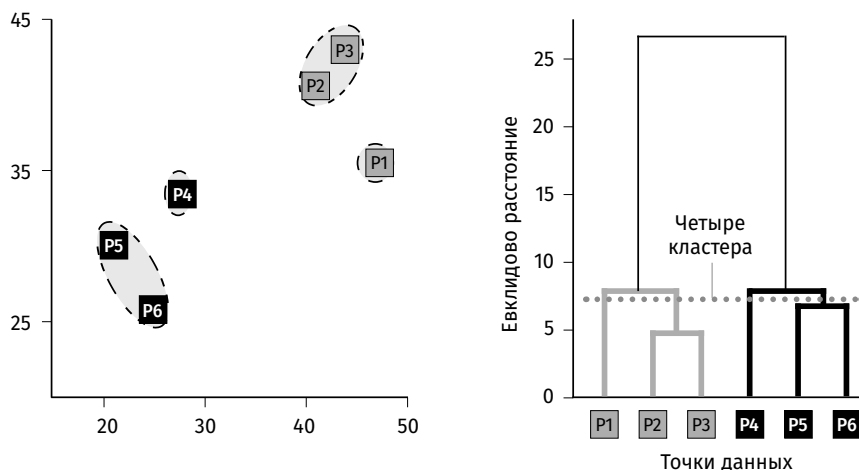


Рис. 6.35. Разделение на четыре кластера

Те вертикальные линии кластеров, которые совпадают с этой пороговой линией или опускаются ниже ее, включаются в наш анализ, а все, что над этой линией, — исключается. В приведенном выше примере P1, P4 и кластеры P2/3 и P5/6 будут включены. Вот каким образом дендрограмма и точки данных (или кластеры) оказываются связаны. Но остается вопрос: как найти оптимальное количество кластеров? Может ли дендрограмма, как и метод локтя, помочь нам выбрать оптимальное число кластеров?

Стандартный метод выполняет поиск вертикальных линий дендрограмм. Он ищет самый длинный вертикальный сегмент, находящийся между *уровнями*, на которых находятся горизонтальные сегменты (это важно — сегмент, который мы ищем, не только не должен прерываться горизонтальными линиями, но и их воображаемыми продолжениями). В нашем случае наибольшее непрерывное вертикальное расстояние показано на рис. 6.36:

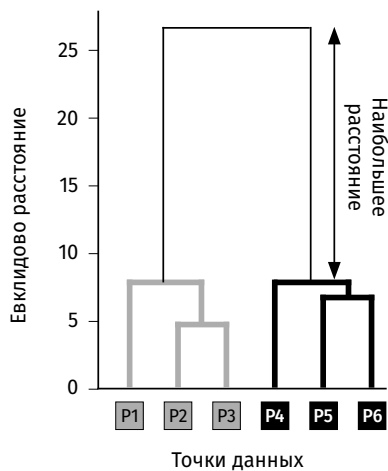


Рис. 6.36. Определение наибольшего вертикального сегмента

После того как вы нашли самую длинную вертикальную линию, установите пороговую линию в точке так, чтобы она пересекала сегмент. Полученное число кластеров оптимально для вашей задачи. В нашем случае это два кластера. Вы, я уверен, согласитесь, что это интуитивно понятно из графика рассеяния (рис. 6.37):

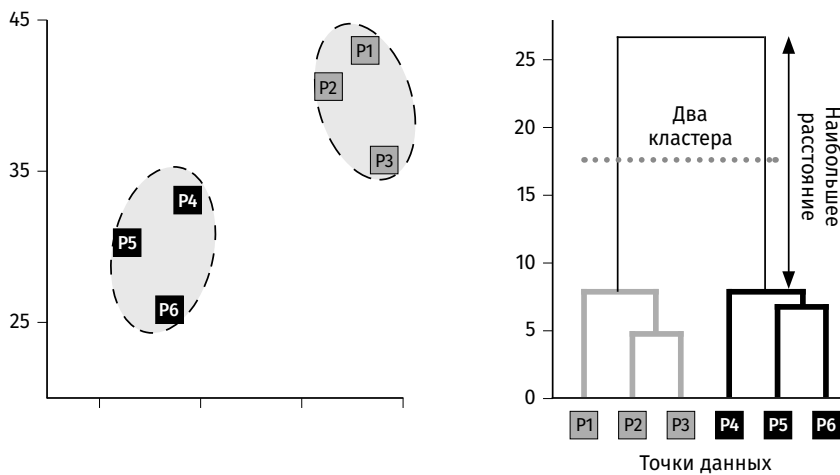


Рис. 6.37. Разделение на два кластера

Самым большим преимуществом использования алгоритма иерархической кластеризации является его дендрограмма. Дендрограмма — практичный визуальный инструмент, который позволяет легко увидеть все потенциальные конфигурации кластера.

Есть еще много алгоритмов как для классификации, так и для кластеризации: те, с которыми я вас познакомил, — только начало. Если вы хотите узнать больше о том, как можно работать с ними в рамках науки о данных, просто посетите [SuperDataScience](#), чтобы найти ряд ресурсов, учебных пособий и курсов.

В следующей главе мы продолжим исследование алгоритмов и рассмотрим один из моих самых любимых типов анализа данных: обучение подкреплением.

Анализ данных (часть II)

07

На протяжении десятилетий средства массовой информации были очарованы роботами — от невинных игрушечных собак до более угрожающих подвидов, замещающих работников физического труда и сотрудников магазинов. Мы значительно продвинулись вперед от созданных в прошлом веке несовершенных автоматов — и во многом благодаря достижениям в сфере обучения с подкреплением. Именно поэтому вторая часть нашего знакомства с анализом данных полностью посвящена алгоритмам из этой области.

Обучение с подкреплением

Обучение с подкреплением — это, по сути, форма машинного обучения, которая опирается на концепцию бихевиоризма при обучении искусственного интеллекта и управлении роботами. Обучение с подкреплением не требует от робота выполнения заданных действий, а позволяет ему исследовать окружающую действительность и обучаться лучшим методам решения задач. Давайте возьмем для примера роботизированную игрушечную собаку. Конечно, можно было бы дать собаке набор пошаговых инструкций, которые скажут ей, что делать, чтобы ходить (передняя правая лапа вперед, левая задняя лапа вперед, левая передняя лапа вперед, правая задняя лапа вперед). Этот метод использовался в более ранних тестах на роботах путем ввода последовательности действий, необходимой для выполнения поставленной задачи. Однако таким образом мы ограничиваем робота этой конкретной последовательностью движений, делая его, ну... только роботом. Но, применяя алгоритм обучения с подкреплением, мы можем заставить собаку-робота самостоятельно научиться ходить.

Использование обучения с подкреплением в случае роботизированных игрушечных собак — самый простой пример того, как этот метод может быть применен. В этой области были сделаны поистине удивительные открытия. Совсем недавно ученые из Лаборатории искусственного интеллекта OpenAI (основанной ведущими технологическими предпринимателями Илоном Маском и Сэмом Алтманом) научили ИИ-ботов выстраивать общую языковую систему, чтобы учиться друг у друга выполнению задач (Recode, 2017). После установки алгоритма обучения с подкреплением боты начали тестировать различные способы связи, чтобы убедиться, что они успешно справятся с поставленной задачей. В ходе этого процесса боты развили общий язык, который основывался на связывании действий, местоположений, объектов и даже самих ботов с абстрактными элементами. Результаты показывают, что ИИ не так уж и отличается от нас: исследователи OpenAI обнаружили, что их боты стремились выполнять задачи более эффективно, развивая свой общий язык таким образом, чтобы он соответствовал проблеме.

Обучение с подкреплением осуществляется путем опробования всех вариантов, доступных машине, а затем отработки оптимальных действий на основе этого индивидуального опыта. В нашем более простом примере, касающемся собаки-робота, ученые, задача которых состоит в том, чтобы собака шагнула вперед и не упала, будут осуществлять обучение с подкреплением, связывая успешное продвижение с наградой, а неудачное продвижение (скажем, падение) с наказанием. В отличие от реальных собак вам не нужно давать реальное поощрение — вы просто отмечаете успешный результат как «1» в своем алгоритме, а неудачный — как «0» или «-1».

Нюансы робототехники заслуживают гораздо больше внимания, чем позволяет объем этой книги. Однако теперь, когда я разжег ваш интерес объяснением того, что такое обучение с подкреплением и как оно может быть применено, мы можем перейти к некоторым конкретным проблемам, которые оно поможет решить, и алгоритмам, способным содействовать этому процессу.

Обучение с подкреплением и поведение человека

Машинное обучение с подкреплением удивительно похоже на процессы усвоения знаний человеком. Один из ярких примеров — то, как ребенок учится ходить. На этой стадии развития малыш действует инстинктивно. Таким образом, мы знаем, что любое поведение управляется системой, хранящейся глубоко внутри мозга, — это нечто бессознательное, нечто, заранее закодированное в нашей ДНК. Как же оно работает?

Когда дети начинают учиться ходить, они часто падают. Обычно падение сопровождается ударом — и нервная система мальчика или девочки посылает сигналы боли в мозг. Таким образом, боль — а это не что иное, как электрический сигнал, посылающийся в мозг, — не существует от нас отдельно. Это чувство, которое создает нервная система, чтобы *тренировать* нас. И так как мозг является частью нашей нервной системы, у нас есть очень интересная установка: когда ребенок падает, одна часть нервной системы дает другой ее части отрицательную обратную связь в виде боли, чтобы заставить понять: подобное действие имеет негативные результаты. В итоге дети узнают, что они не должны повторять действие, приведшее к падению. Сногшибательно, если задуматься.

В то же время, если ребенку удастся сделать шаг вперед — скажем, чтобы поймать кошку за хвост или дотянуться до конфеты на краю стола, — его нервная система пошлет положительные сигналы в мозг и малыш будет вознагражден. Повторяя эти действия, он научится ходить.

Поразительно: мы создаем алгоритмы обучения с подкреплением в области ИИ и робототехники, в то время как наша собственная нервная система является самым впечатляющим из всех алгоритмов обучения с подкреплением.

Задача о «многоруком бандите»

«Многорукий бандит» может напомнить персонаж, которого вы ожидаете встретить в «Игре престолов», но это просто обозначение общей

задачи, связанной с обучением с подкреплением. Популярные решения задачи о «многоруком бандите» — алгоритм верхней доверительной границы и выборка Томпсона, каждый из которых мы здесь рассмотрим.

Так откуда же такое образное название? Задача о «многоруком бандите» была впервые сформулирована с отсылкой к самым первым игровым автоматам в казино. Вместо привычных нам сегодня игровых автоматов с кнопочным управлением у более ранних моделей имелся рычаг сбоку, который игрок тянул, чтобы вращать (обычно) три барабана. Из-за этого единственного рычага — а также потому, что такие машины исключительно быстро вынуждали игроков оставить все свои деньги в казино, — автоматы стали называть «однорукими бандитами».

Первые модели игровых автоматов давали игрокам почти 50%-ный шанс выиграть или проиграть. Через некоторое время казино переоснастили свои машины, чтобы значительно уменьшить шансы игроков. Будет ли когда-нибудь возможно победить систему, если преимущества не на стороне игроков? В связи с ответом на этот вопрос появилась задача о «многоруком бандите». Обычно в казино есть несколько игровых автоматов (таким образом, «бандит» оказывается многоруким, а не одноруким). Если вероятность выигрыша различна для каждого из этих автоматов и мы не знаем, у какого из них она выше, то как нам играть с выбранным количеством игровых автоматов в определенном порядке, чтобы максимизировать выигрыш?

Задачу о «многоруком бандите» можно рассматривать гораздо шире: с ее помощью можно провести наиболее эффективную рекламную кампанию (алгоритмы, которые мы будем изучать, отличны от пресловутого А/В-теста, случайного эксперимента, в котором два варианта — А и В — противопоставляются друг другу, чтобы определить оптимальный), наиболее эффективным образом выделить ресурсы на исследовательские проекты или помочь усовершенствованию эксплуатационных функций роботов.

Верхняя доверительная граница и А/В-тестирование

Столь высокая эффективность метода обучения с подкреплением обусловлена тем, что он использует варианты, доступные благодаря подходу, который ориентируется на данные. С другими способами тестирования — такими, как А/В-тестирование, часто применяемое в маркетинге, — решение может быть принято только после того, как все варианты изучены равное количество раз, и тогда, когда у нас есть достаточно большая выборка, на основе которой мы можем делать уверенные выводы. На изучение каждого варианта таким единым образом тратится много времени и денег, в то время как другие алгоритмы, в частности верхняя доверительная граница, могут подойти к поиску нашего оптимального результата путем *динамического* проведения тестов, включающих в себя как исследование (случайный выбор), так и использование (выбор на основе предварительных знаний). Мы рассмотрим и то и другое более подробно далее в этой главе. Такой подход призван не только максимально быстро найти оптимальный вариант, но и максимизировать вашу прибыль в процессе работы. В принципе, алгоритм верхней доверительной границы выглядит предпочтительнее тестирования А/В.

Тестирование задачи

А пока давайте обратимся к ярким огням Вегаса. Я не одобряю и не поощряю азартные игры; этот пример — просто отличный способ показать ход решения задачи о «многооруком бандите».

Итак, при полном параде мы вошли в казино Caesars Palace, и перед нами пять игровых автоматов. В каком порядке и сколько раз мы должны играть на них, чтобы максимизировать выигрыш? Сначала допустим, что для каждой машины заранее задано распределение результатов (проигрышей и выигрышей). После того как мы потянем за рычаг (или нажмем кнопку), выбранный нами игровой автомат случайным образом выберет результат (выигрыш или проигрыш) согласно распределению — скажем, если вы поставите 50 центов, то либо получите обратно \$0 (проигрыш), либо \$1 (выигрыш).

Все, что нам нужно знать, — это распределение вероятностей выигрыша на каждом игровом автомате, чтобы играть исключительно на том, который дает наиболее благоприятные шансы на победу*. Легко.

Но вот в чем проблема: мы *не знаем* этих распределений и Caesars Palace вряд ли предоставит эту информацию паре оптимистичных аналитиков данных!

Ставки для решения этой реальной проблемы в казино высоки. Мы должны потратить наши деньги на проведение экспериментов, и чем дольше будем искать решение, тем больше денег потратим. По этой причине мы должны найти нужный результат как можно быстрее, чтобы сократить наши потери.

Для поддержания эффективности следует учитывать два фактора — исследование и эксплуатацию — и применять их в тандеме: исследование означает поиск лучшей машины, а эксплуатация — применение знаний, которые у нас *уже* есть о каждой из машин, чтобы делать ставки. Дело в том, что без предварительной разведки у нас не будет данных для применения, а без применения мы станем зарабатывать меньше денег, чем могли бы в случае опоры на собранную информацию.

Прежде чем мы начнем с азов верхней доверительной границы и выборки Томпсона, давайте изучим решение задачи о «многоруком бандите». Во-первых, будем считать, что имеем дело с пятью игровыми автоматами. Если иллюстрировать произвольный набор распределений для них, диаграмма могла бы выглядеть примерно так, как на рис. 7.1.

Этот график иллюстрирует вероятность проигрыша (\$0) и выигрыша (\$1) от каждого из пяти игровых автоматов. Например, если вы вставляете 50 центов в машину D3, есть 90%-ная вероятность того, что вы получите \$0, и 10%-ная — того, что разбогатеете на \$1**.

* Предположим для этого примера, что существует по крайней мере один игровой автомат, который в конечном итоге позволяет игроку выиграть больше, чем проиграть. Это не противоречит сказанному ранее: в совокупности игровые автоматы все еще могут быть настроены таким образом, чтобы в более чем в половине случаев казино выиграло в целом.

** Такое распределение вероятностей, которое описывает результаты «да»/«нет», называется распределением Бернулли.

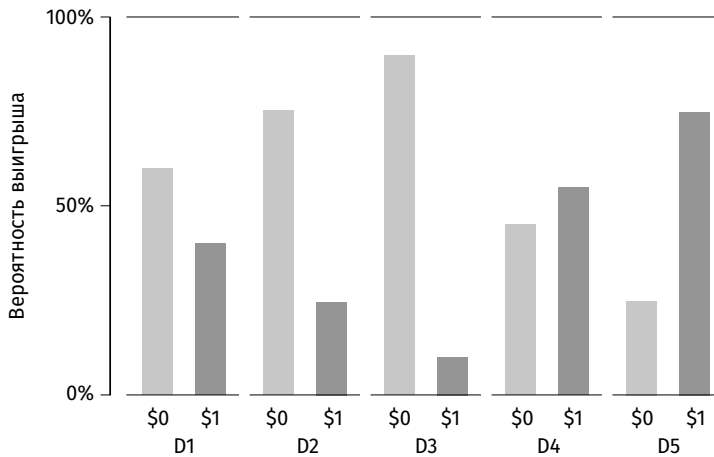


Рис. 7.1. Распределение вероятностей выигрыша при игре на «многоруких бандитах»

Из графика становится очевидно, что игровой автомат D5 в среднем даст наилучший результат, потому что у него наивысшая вероятность (75%) выигрыша.

Чтобы упростить, мы можем рассчитать и построить график ожидаемого выигрыша для каждой из машин по формуле:

$$E(X) = (p \times \text{результат 1}) + (q \times \text{результат 2}),$$

где $E(X)$ — ожидаемый выигрыш, p — вероятность выигрыша, q — вероятность проигрыша, а результат 1 и результат 2 — суммы, полученные в случае выигрыша и проигрыша соответственно.

Например, для машины D2 ожидаемый выигрыш будет рассчитываться как:

$$(25\% \times \$1) + (75\% \times \$0) = \$0,25.$$

Самый простой способ прикинуть размер ожидаемого выигрыша таков: если вы долго играете на машине D2, а затем усредняете свои результаты по количеству игр, то теоретически среднее значение, которое вы вычисляете, должно быть равно \$0,25. В сущности, ожидаемый выигрыш — это ваш теоретический средний выигрыш.

Чем выше ожидаемый выигрыш машины, тем больше у вас шансов выиграть на ней, и наоборот. Ожидаемый выигрыш для всех пяти

машин указан на рис. 7.2. В качестве быстрого упражнения проверьте правильность вычислений.

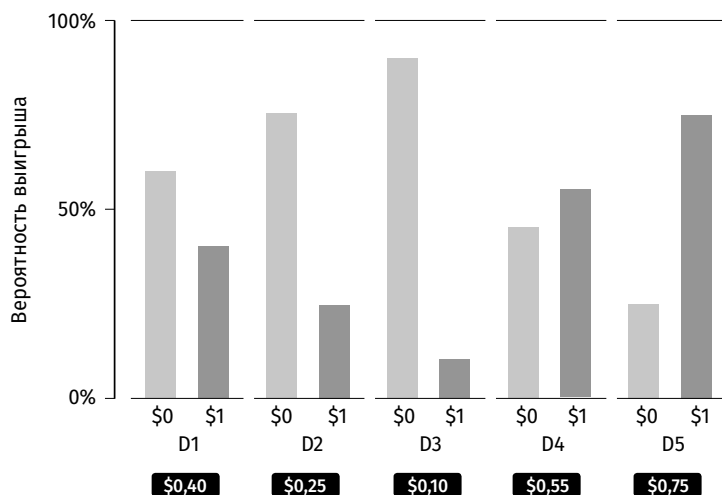


Рис. 7.2. Распределения «многоруких бандитов» с ожидаемым результатом

Опять же, эти распределения и ожидаемый выигрыш представлены здесь в качестве примера. Если бы мы заранее знали ожидаемый выигрыш, то пошли бы прямо к пятой машине и постоянно делали ставки на ней, так как в среднем получили бы наилучшие результаты. Реальность такова, что мы сможем узнать этот ожидаемый выигрыш (или приблизиться к его нахождению) только в процессе игры.

Это одна из удивительных сторон науки о данных, и именно поэтому задача о «многоруком бандите» — одна из моих любимых. Разработка лучших решений и построение массива данных для них путем проб и ошибок, исследования и использования напоминают мне о тех методах, которые Алан Тьюринг и команда Блетчли-парка использовали для нахождения кода шифра. В мире бизнеса эти методы часто используются в онлайн-рекламе и маркетинге для тестирования реакции потребителей на различные варианты плакатов, созданных для рекламирования одного продукта.

Машины тоже рискуют

Наша цель при решении задачи о «многоруком бандите» — найти и использовать лучший игровой автомат, а также потратить наименьшее количество времени (и денег) на изучение всех вариантов. Поскольку мы не знаем заранее, какой из «бандитов» лучший, мы получаем то, что известно в математике как риск. По существу, риск — это то, что происходит в ходе эксперимента, когда мы используем неоптимальные методы. В случае с нашими игровыми автоматами в Вегасе игра на неоптимальных машинах сопряжена с риском, определяющимся количественно как разница между оптимальным и неоптимальным результатами. С практической точки зрения риск представляет собой альтернативные издержки изучения игровых автоматов, которых нам удалось бы избежать только с помощью магии. Чем дольше вы исследуете неоптимальные варианты, тем выше риск.

Риск перестает накапливаться, когда мы находим оптимальную машину. Тем не менее существует опасность того, что если мы не будем исследовать все пять «бандитов» достаточно долго, то неоптимальная машина может показаться оптимальной — и из-за нашей спешки риск будет нарастать. Поэтому важно потратить время на изучение, прежде чем переходить к каким-либо выводам. На примере наших пяти распределений мы видим, что D4 также достаточно хороший игровой автомат для получения прибыли. Однако он не является оптимальным. Если бы мы не знали об этих распределениях вероятностей выигрыша и сначала должны были бы исследовать машины, не полностью изучив все пять вариантов, ожидаемый выигрыш от D4, маскирующейся под оптимальный автомат, мог бы ввести нас в заблуждение. Конечно, мы бы продолжали получать деньги, и определенно больше, чем могли бы заработать на машинах 1, 2 и 3, но это был бы не лучший результат из возможных.

Теперь, когда мы узнали, что такое задача о «многоруком бандите» и как ее можно использовать в работе, попробуем ответить на вопрос: какие алгоритмы мы можем применить к проектам, требующим таких инструментов? Два наиболее распространенных метода, по крайней мере в мире бизнеса, — верхняя доверительная граница и выборка Томпсона. В обоих алгоритмах проводится последовательная проверка

различных вариантов и ведется регистрация результатов в лог-файлы, содержащие лучшие и худшие распределения. В дальнейшем мы изучим нюансы каждого алгоритма, а также преимущества и недостатки их использования.

Верхняя доверительная граница

Предупреждаю: мы продолжим рассматривать пример игровых автоматов, но его следует воспринимать только как гипотетический. Проведение подобного эксперимента в Лас-Вегасе в лучшем случае сделает вас не слишком популярным, а в худшем — приведет в тюрьму. Применяя к нашей задаче алгоритм верхней доверительной границы (ВДГ), мы определим, у какой машины лучший ожидаемый выигрыш, — а это высветит ключевое различие между алгоритмами, рассмотренными в предыдущей главе, и теми, о которых говорится здесь. В наших предыдущих примерах мы, как правило, использовали массивы данных с собранными независимыми и зависимыми переменными. Однако в обучении с подкреплением все по-другому. Мы начинаем вообще не с данных. Мы должны экспериментировать, наблюдать и менять нашу стратегию на основе предыдущих действий.

Когда вы выигрываете, алгоритм верхней доверительной границы фиксирует в своем массиве данных получение вами выигрыша как 1. Потери будут записаны как 0. Для каждой игры ВДГ добавит результат в свой массив данных. Так алгоритм обучается с помощью исследования и одновременно разрабатывает стратегию, которая позволит избежать случайного выбора машин*. Ход выполнения алгоритма от одного автомата к другому будет зависеть от результатов каждого предыдущего раунда, а динамическая стратегия повышает точность при сборе дополнительной информации. Например, машина, выбранная алгоритмом в 281-м раунде нашего теста, будет выбрана на основе всех данных, собранных за предыдущие 280 раундов.

* Естественно, на первом этапе неизбежен некоторый уровень случайности, так как у нас не будет никаких данных на этот момент, но должны же мы с чего-то начать, в конце концов!

Построение алгоритма верхней доверительной границы

Прежде всего предварительно: мы начинаем с определенного количества «рук» — в нашем случае это соответствует пяти игровым автоматам в Вегасе. Выбор одного из пяти автоматов представляет собой «раунд», или «игру» (я буду использовать оба термина). Каждый раз, когда мы опускаем деньги в машину и тянем за рычаг, мы завершаем раунд. Как только раунд завершен, мы либо будем вознаграждены, либо потеряем свои деньги. Эта информация запишется алгоритмом верхней доверительной границы как 1 для выигрыша или 0 для потери. Как мы обсуждали, наша цель — найти машину с наивысшим ожидаемым выигрышем. Давайте покажем результаты, которые мы ищем, на диаграмме (рис. 7.3)*.

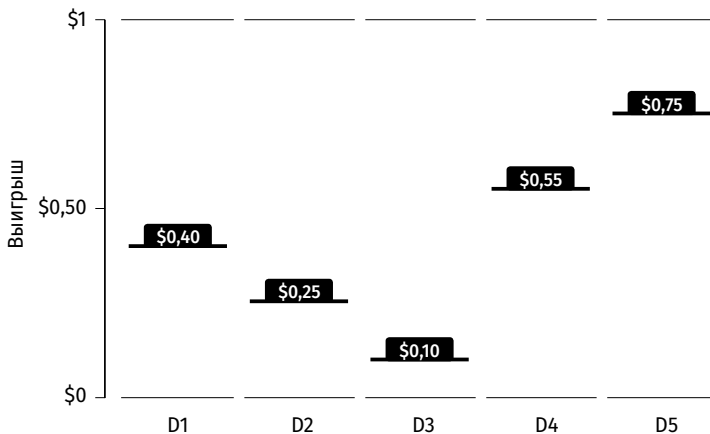


Рис. 7.3. Вертикальная ось ожидаемого выигрыша

1. Выберите начальную точку

Верхняя доверительная граница задает универсальную стартовую точку для всех «рук», и она основана на предположении, что все «руки» будут давать одинаковые результаты. В случае наших игровых автоматов алгоритм будет устанавливать в качестве стартового значения среднюю точку между выигрышем и потерей: 0,5 (пунктирная линия). Затем эта линия смещается для каждой машины, по мере того как будут сыграны успешные игры (рис. 7.4).

* Помните: это упражнение поможет нам понять, как работает алгоритм. В реальном сценарии мы не знали бы об ожидаемом выигрыше.

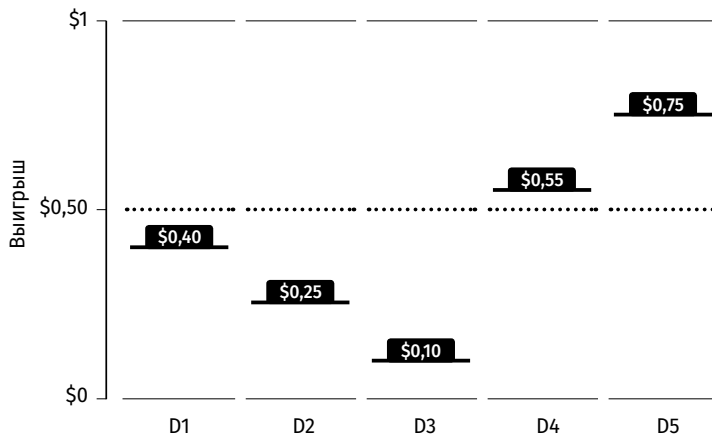


Рис. 7.4. Стартовая линия

2. Установите начальную доверительную границу

Алгоритм установит верхние и нижние доверительные границы, которые охватывают все возможные ожидаемые выигрыши. В нашем примере (рис. 7.5) это означало бы, что верхняя доверительная граница будет соответствовать 1 (выигрыш), а нижняя доверительная граница — 0 (проигрыш). Мы можем быть уверены, что верхняя доверительная граница должна закончиться здесь, потому что наши игры не могут привести к чему-то большему, чем победа. Та же логика

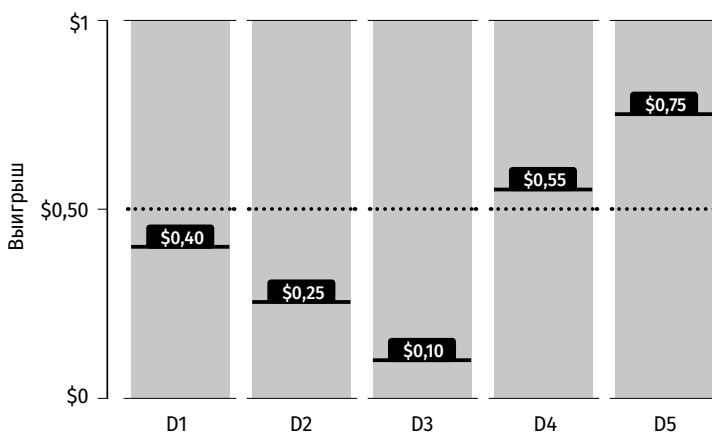


Рис. 7.5. Начало доверительной границы

относится к тому, что наши игры не могут привести к чему-то меньшему, чем проигрыш.

Важно понять предназначение границы доверия. В реальной ситуации мы не знали бы точно, где находится ожидаемая прибыль. В начале первой игры нам вообще было бы это неизвестно. То, что алгоритм «видит» в начале первого раунда, выглядит примерно так, как показано на рис. 7.6:

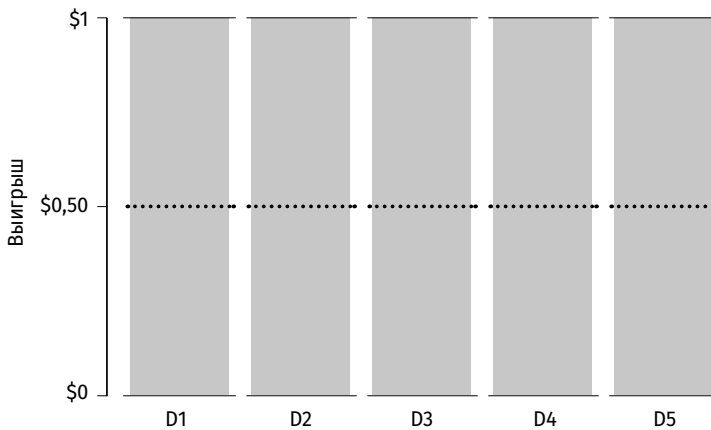


Рис. 7.6. Реальный сценарий

Доверительные границы устанавливают для ожидаемого выигрыша определенные рамки. Например, ожидаемый выигрыш машины D2 может быть меньше или больше \$0,5 — точнее сказать невозможно. Однако серый прямоугольник доверительных границ построен таким образом, что мы знаем: ожидаемый выигрыш от D2 должен быть где-то внутри него. На данном этапе это очевидно, так как выигрыш не может быть меньше \$0 или более \$1. Однако, как вы увидите далее, когда игры будут воспроизводиться на пяти машинах, алгоритм продолжит перемещать и изменять размеры этих доверительных границ так, чтобы они отображали ожидание выигрыша. Это ключ к алгоритму ВДГ*.

* Чтобы быть полностью корректными — мы никогда не можем быть на 100% уверены в точном диапазоне ожидаемого выигрыша. Он может быть где угодно! И именно поэтому границы рассчитываются таким образом, что ожидаемая доходность попадает в них с определенным уровнем уверенности (например, 95%). Это также объясняет, почему они называются доверительными границами.

3. Сыграйте пробные раунды

Первые несколько раундов будут пробными — благодаря им у нас соберутся изначальные данные, которые затем мы используем для информированного принятия решений в более поздних раундах. На этом этапе придется сыграть несколько раз на каждой из машин, чтобы можно было перенастроить доверительные границы для каждого из распределений.

Предположим, что мы сыграли на машине D3 10 раз, что создало следующую последовательность побед и поражений: 1, 0, 0, 0, 0, 0, 1, 0, 1, 0. Значит, *наблюдаемый* средний выигрыш окажется следующим:

$$(1 + 0 + 0 + 0 + 0 + 0 + 1 + 0 + 1 + 0)/10 = \$0,30.$$

Истинный выигрыш vs наблюдаемый выигрыш

Под истинным ожидаемым выигрышем (или просто «ожидаемым выигрышем») мы понимаем ожидаемый *теоретический* выигрыш в соответствии с распределением, запрограммированным на каждой машине, как показано на рис. 7.2. Его не следует путать с *наблюдаемыми* средними выигрышами, которые мы будем рассчитывать на основе данных, полученных при взаимодействии с машинами. Первый — теоретический; второй — эмпирический.

Теперь, когда игры сыграны, алгоритм сдвинет пунктирную линию D3 до \$0,30 (рис. 7.7).

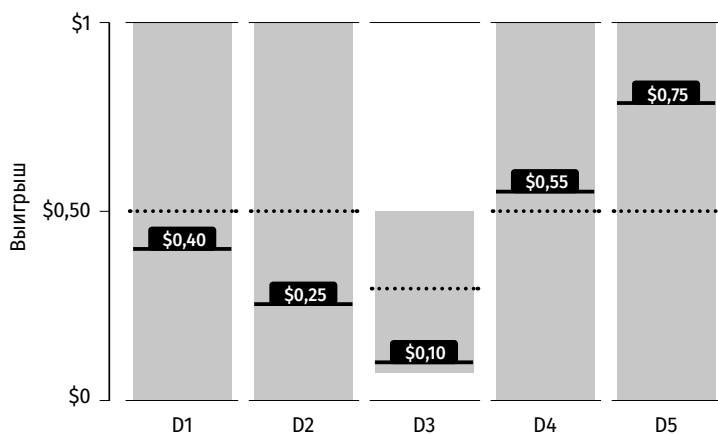


Рис. 7.7. Игровой автомат D3 после пробных раундов

Но что происходит с серым прямоугольником, который представляет наши верхние и нижние доверительные границы? Он будет *сжиматься* каждый раз, когда мы играем раунд. Это связано с тем, что чем больше раундов мы играем, тем более точным будет наблюдаемое среднее значение и, следовательно, более узкими доверительные границы. Тогда размер этого прямоугольника окажется обратно пропорционален количеству раундов, которые мы сыграли на данной машине. Чем меньше прямоугольник, тем более верно то, что мы приближаемся к истинному ожидаемому выигрышу данной машины. Это прямое следствие закона больших чисел.

Закон больших чисел

Закон больших чисел (ЗБЧ), безусловно, мой любимый закон в математике: он интуитивен и применим к реальной жизни. Мы не будем вдаваться в подробности, но проиллюстрируем основные концепции ЗБЧ.

Скажем, вы подбрасываете монету. Вероятность выпадения орлов или решек составляет 50/50, но это не гарантирует, что распределение результатов игры всегда будет равным. Иногда выпадет больше орлов, а иногда — больше решек. Например, если вы будете бросать монету 10 раз, вполне вероятно получить семь орлов и три решки. Итоговое распределение орлов/решек при таком исходе составляет 70/30%.

Но если вы бросите монету 100 раз, какова вероятность получить 70 орлов и 30 решек? Она намного ниже. Подумайте об этом — если вы получите 70 орлов из 100 подброшенных монет, разве вы не подумаете, что дело нечисто? Было бы гораздо более правдоподобно, если бы у вас оказалось 59 орлов и 41 решка, что составило бы 59/41%.

Далее, что может произойти, если вы подбросите монету 1000 раз? Вероятность получить 700 орлов и 300 решек в честном споре будет почти равна нулю. Опять же, не вдаваясь в математику, — как бы вы себя чувствовали, если бы вам выпало 700 орлов из 1000 подброшенных монет? Гораздо более реалистичным был бы результат, если бы вы получили (скажем) 485 орлов и 515 решек с распределением 48,5/51,5%.

Таким образом, мы можем видеть, что по мере увеличения количества подбрасываний монет *наблюдаемое* распределение приближается к *истинному* ожидаемому распределению 50/50. Это закон больших чисел: по мере роста размера выборки наблюдаемое среднее всегда будет приближаться к истинному ожидаемому результату.

Предположим, что мы сыграли десять игр на каждом из наших игровых автоматов. Результаты отображены на рис. 7.8.

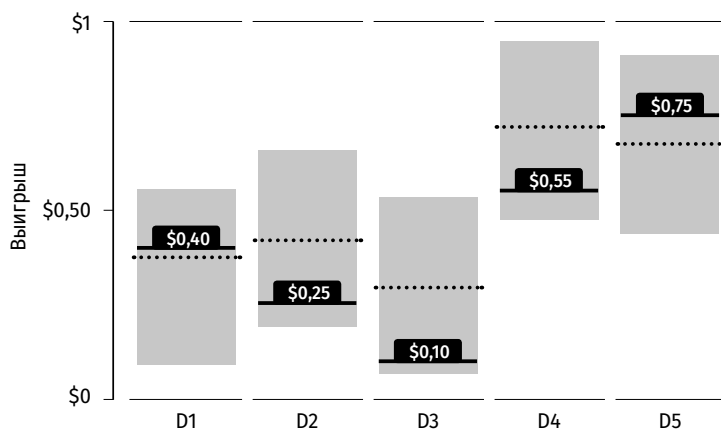


Рис. 7.8. Все игровые автоматы после пробных раундов

4. Определите оптимальный автомат и играйте на нем

Как только будет достаточно данных по всем игровым автоматам, алгоритм верхней доверительной границы начнет разворачивать анализ в направлении машин, имеющих наивысшую верхнюю доверительную границу, при этом неоптимальные «руки» не будут учитываться. Это интуитивно понятно: поскольку истинный ожидаемый выигрыш машины может быть каким угодно в пределах ее доверительных границ, алгоритм предполагает, что оптимальным автоматом будет тот, у которого самая высокая верхняя доверительная граница (отсюда и название). В нашем примере оптимальной машиной представляется D4.

Однако, глядя на истинные ожидаемые выигрыши, вы увидите, что D4 явно неоптимальна. Не волнуйтесь, границы защитят нас от выбора неоптимальной машины в долгосрочной перспективе. Если мы играем на неоптимальном автомате достаточно долго, его наблюдаемое среднее значение приблизится к ожидаемому выигрышу и прямоугольник на диаграмме в конечном итоге будет сведен к точке, где алгоритм сочтет оптимальной другую машину. Это связано с тем, что другой игровой автомат, на котором не играли так часто, будет иметь гораздо более широкие доверительные границы. В нашем случае, когда наблюдаемое среднее значение D4 приближается к ожидаемому выигрышу для этой машины, а ее доверительные границы достаточно узки, алгоритм переключается на автомат D5 (рис. 7.9).

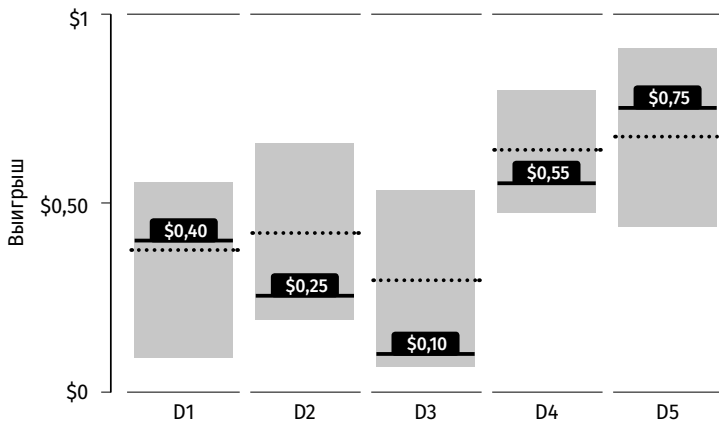


Рис. 7.9. Сдвиг и сужение доверительного предела D4

Пока мы используем машину D5, ее доверительные границы также будут сужаться, и алгоритм может даже вернуться к D4 на некоторое время. Однако это только до тех пор, пока доверительные границы D4 не станут достаточно узкими. В конечном счете верхние доверительные границы оптимальной машины будут оставаться выше верхних доверительных границ всех других машин («рук»).

Алгоритм верхней доверительной границы подходит для:

- поиска наиболее эффективных рекламных кампаний;
- управления большим числом финансовых проектов.

ВДГ не единственный алгоритм, который может решить проблему «многорукого бандита». Далее мы рассмотрим, как можно применить выборку Томпсона, — подумайте о том, когда этот алгоритм может оказаться предпочтительнее ВДГ.

Выборка Томпсона

Прежде чем продолжить, хочу отметить одну важную вещь. Понимание алгоритма ВДГ поможет нам уяснить методы, которые мы применяем, поэтому, если вы читаете разделы выборочно, я рекомендую полностью прочитать все, что относится к обучению с подкреплением. При этом уделите особое внимание вопросам, связанным с задачей о «многоруком бандите», чтобы как можно лучше усвоить выборку Томпсона*. (Имейте в виду: выборку Томпсона понять труднее, чем алгоритм верхней доверительной границы. Если хотите, можете пропустить и изучить позже посвященный ей раздел.)

Помните, что задачу о «многоруком бандите» мы решаем для того, чтобы наиболее эффективно исследовать и использовать наши варианты и тем самым максимально увеличить выигрыш. В данном примере (рис. 7.10) мы облегчим задачу и возьмем три игровых автомата,

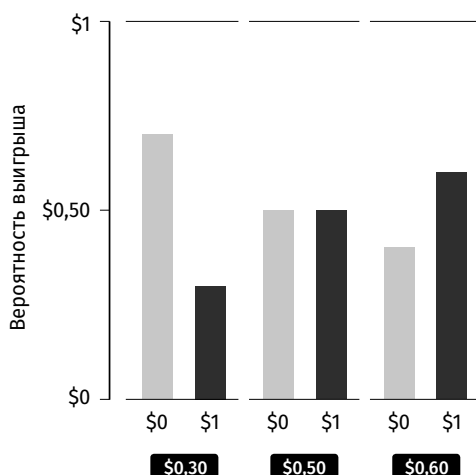


Рис. 7.10. Выборка Томпсона для распределения «многоруких бандитов»

* Названа в честь разработчика Уильяма Томпсона.

поскольку выборка Томпсона чуть сложнее по сравнению с верхней доверительной границей.

Эта диаграмма аналогична той, что мы видим на рис. 7.2 в разделе о верхней доверительной границе. Здесь ожидаемый выигрыш от машины M1 рассчитывается как

$$(0,7 \times \$0) + (0,3 \times \$1) = \$0,30.$$

Ожидаемые выигрыши для машин M2 и M3 рассчитываются с использованием того же метода.

Построение выборки Томпсона

Как и в случае с алгоритмом ВДГ, давайте начнем с построения графиков ожидаемых выигрышей от каждой из наших трех машин.

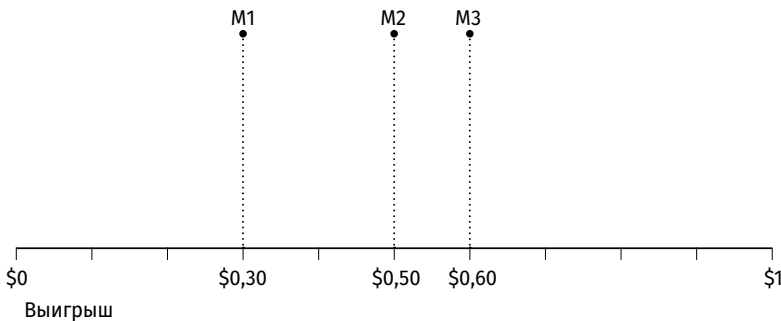


Рис. 7.11. Выборка Томпсона для ожидаемого выигрыша

На этом графике (рис. 7.11) на оси x показан ожидаемый выигрыш, а вертикальные пунктирные линии представляют собой наши автоматы M1, M2 и M3*. Как и в случае с алгоритмом верхней доверительной границы, эти строки отражают *истинный* ожидаемый выигрыш от машин. Хотя распределение вероятностей выигрыша для каждой из машин и отображено на графике, нам оно неизвестно (если только мы не связаны с владельцем казино). Поэтому в реальной жизни эти ожидаемые значения были бы нам неизвестны; наша цель — найти

* Это немного отличается от графика для алгоритма ВДГ, где ожидаемый выигрыш показан на оси y .

их. Мы показываем их здесь, чтобы продемонстрировать предсказательную способность алгоритма выборки Томпсона.

1. Сыграйте пробные раунды

Прежде чем мы сможем проверить наши данные, их нужно сначала собрать. Это означает, что вы играете несколько раундов, чтобы иметь возможность оценивать игровые автоматы. Скажем, мы сыграли на трех автоматах по 12 раз, и машина М3 дала следующую последовательность побед и поражений: 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1. Это значит, что наш средний выигрыш составляет \$0,67. Но алгоритм выборки Томпсона не так прост: он знает, что это всего лишь *наблюдаемый* средний выигрыш и что истинный ожидаемый выигрыш необязательно составляет \$0,67. Пока размер выборки мал и мы можем только сказать, что истинный ожидаемый выигрыш равен примерно этой сумме. Для решения задачи алгоритм выборки Томпсона построит кривую распределения вероятности, чтобы оценить, где может быть истинный ожидаемый выигрыш (рис. 7.12).

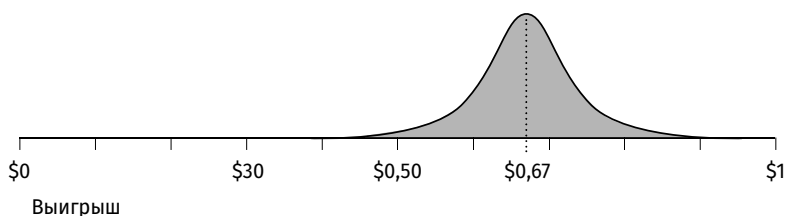


Рис. 7.12. Распределение вероятностей выигрыша для игрового автомата М3 после пробных раундов

Мы видим, что распределение сосредоточено вокруг \$0,67. Это означает, что на основе данных, которыми располагает алгоритм, он оценивает, что истинный ожидаемый выигрыш может быть либо равен \$0,67, либо близок к этому значению. Чем дальше от \$0,67, тем ниже вероятность того, что мы имеем дело с истинным ожидаемым выигрышем. Это разумная оценка, потому что если бы ожидаемый выигрыш составлял, например, \$0,1, то едва ли в наших пробных играх мы бы восемь раз выиграли и только четыре — проиграли; вместо этого у нас было бы намного меньше побед.

Если мы теперь добавим на график истинное ожидание выигрыша (которого алгоритм не знает), то увидим, что истинный ожидаемый выигрыш от автомата М3 довольно близок к центру кривой распределения (рис. 7.13).

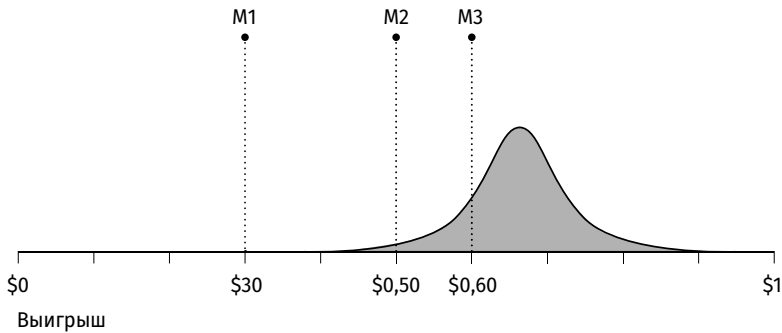


Рис. 7.13. Распределение вероятностей выигрыша для игрового автомата М3 с истинными ожидаемыми выигрышами

Обратите внимание: все, что мы сделали до сих пор, очень похоже на построение алгоритма ВДГ. Только вместо распределений были доверительные границы (прямоугольники), в которые должен попадать истинный ожидаемый выигрыш. Теперь давайте построим две другие кривые распределения после начальных 12 пробных игр (рис. 7.14).

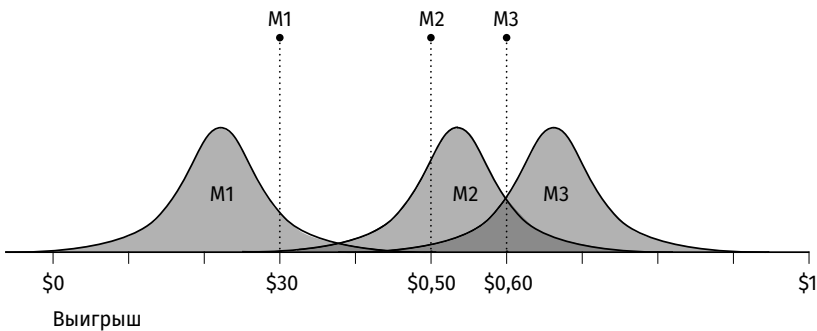


Рис. 7.14. Распределение вероятностей выигрыша для всех трех игровых автоматов после пробных раундов

2. Выберите случайные значения из распределений, чтобы получить задачу о «многоруком бандите»

Вот где начинается веселье. В начале нового раунда выборка Томпсона сначала будет отображать случайное значение из построенного для каждой машины распределения и использовать эти значения для создания своего собственного гипотетического «взгляда на мир». Этот этап очень важен, поскольку мы не знаем, где находится истинный ожидаемый выигрыш. Однако благодаря полученным распределениям вероятностей знаем, где эти выигрыши могут быть. Вот почему мы получаем величины всех распределений и делаем предположение, что они равны или что мы знаем истинное ожидание. В некотором смысле мы создали мнимую вселенную, и теперь нам нужно решить задачу внутри нее. Учитывая характер кривой распределения, вполне вероятно, что алгоритм возьмет точку данных из области, где находится самая высокая кривая*. Однако также возможно, что точки берутся из концов хвоста кривой, как мы можем видеть на рис. 7.15.

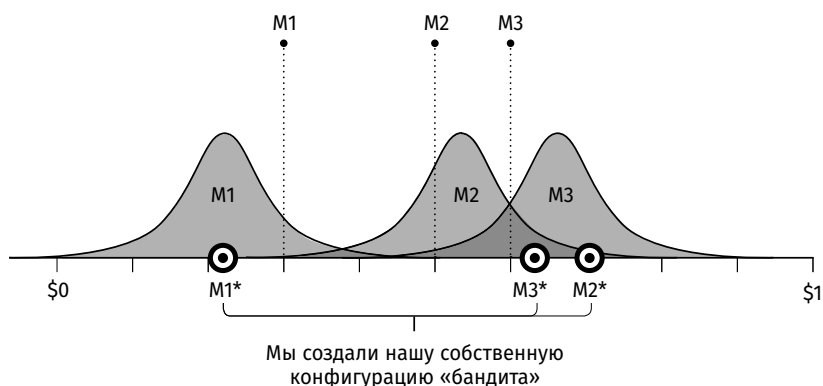


Рис. 7.15. Случайно сгенерированная конфигурация «бандита»

Вышеуказанные три точки данных (M1*, M2* и M3*) представляют гипотетическую конфигурацию алгоритма ожидаемых выигрышей

* Представьте себе область под кривой распределения как пространство, в котором может быть нарисована точка данных. Логически, вы в состоянии увидеть, что существует больше возможностей для точки данных там, где самая высокая кривая, потому что она занимает больше места на графике, чем хвосты распределения.

для каждой машины. Поскольку мы предполагаем, что это правильный взгляд на мир (то есть что $M1^*$, $M2^*$ и $M3^*$ являются истинными ожидаемыми выигрышами), решение задачи о «многоруком бандите» теперь становится тривиальным: точка данных $M2^*$ — самая дальняя на оси x , поэтому автомат $M2$ даст нам лучший результат в этом раунде.

Вероятностное и детерминированное обучение с подкреплением

В главе 6 мы узнали о вероятностном и детерминированном подходах. Они используются в аналитике довольно часто, и было бы полезно напомнить концептуальные различия между ними.

Выборка Томпсона вероятностна, тогда как алгоритм верхней доверительной границы детерминирован — и легко понять почему. Оба подхода похожи тем, что во время игры они приближают нас к значению истинного ожидаемого выигрыша. ВДГ делает это через доверительные границы, тогда как выборка Томпсона создает распределения. Однако ВДГ работает по жестким правилам; когда нам нужно выбрать автомат для игры, мы просто берем машину с наивысшей верхней доверительной границей. При выборке Томпсона вместо (детерминированного) выбора «бандита» для игры в начале каждого раунда мы извлекаем значения из распределения вероятностей и основываем выбор машины на этих значениях.

Если бы мы дважды применили алгоритм верхней границы доверия к одной и той же проблеме, оба раза мы получили бы одинаковый результат после идентичной последовательности итераций. Если, однако, мы дважды применим выборку Томпсона к одной и той же задаче, то, вероятно, в каждом случае получим один и тот же результат (то есть выберем оптимальную машину), но способ, которым были сыграны раунды, был бы совершенно другим, потому что мы каждый раз произвольно генерируем гипотетических «бандитов». Таково ключевое различие между детерминированным и вероятностным подходами.

3. Играйте на «оптимальной» машине

Основываясь на нашей гипотетической конфигурации ($M1^*$, $M2^*$ и $M3^*$), мы можем теперь сыграть раунд на «оптимальном» игровом автомате. Затем будут получены данные (либо выигрыш, либо проигрыш), обновляющие кривую распределения. Предположим, что, когда мы играли на машине M2, итогом был проигрыш (ноль). Этот ноль будет добавлен к ряду результатов, которые мы получаем от этой машины, и он обновит распределение для M2 (рис. 7.16).

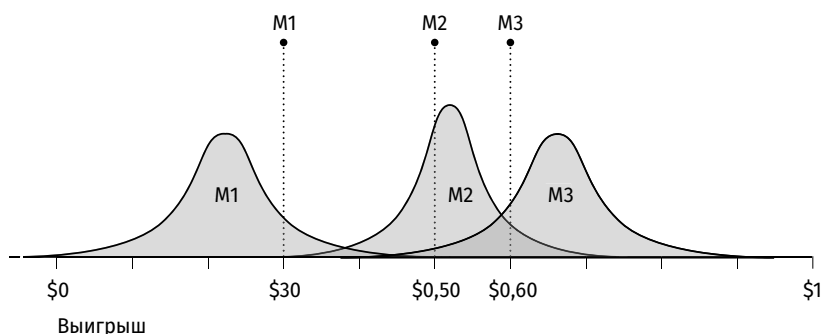


Рис. 7.16. Обновленное распределение вероятностей выигрыша для игрового автомата M2

Нулевой результат для этого раунда уменьшает наблюдаемый средний выигрыш от автомата M2, поэтому кривая распределения смещается влево*. Мы также видим, что кривая сузилась — она больше сосредоточена вокруг центра. Это связано с тем, что игра в дополнительном раунде увеличила размер выборки для этой машины, и, как мы теперь знаем из закона больших чисел, больший размер выборки означает, что мы можем быть более уверены в том, что близки к нахождению истинного ожидаемого выигрыша.

Стоит отметить, что с этой новой кривой распределения менее вероятно, что машина M2 будет иметь наивысший мнимый ожидаемый выигрыш $M2^*$, когда мы станем ставить задачу о «многоруком

* Таково обучение с подкреплением в действии. Алгоритм сделал неверный выбор — и наказывается неудачей (нолем). Распределение смещается влево, чтобы алгоритм помнил, что сделал потенциально плохой выбор и должен попытаться избежать его в будущем.

бандите» для следующего раунда. Это связано с тем, что ее кривая распределения теперь больше сдвинута влево и сузилась, поэтому вероятно, что значение, выбранное из распределения автомата M3, будет больше, чем значение, выбранное из распределения автомата M2. Это относительное расположение кривых распределения M2 и M3 согласуется с реальным состоянием вещей: истинный ожидаемый выигрыш M3 больше, чем M2.

4. Продолжайте играть раунды, чтобы уточнить построенные кривые распределения

Теперь мы можем сыграть дополнительные раунды. Каждый раз, когда мы будем играть, алгоритм еще раз выберет три точки данных для нашей конфигурации и выявит лучшую из них (самую дальнюю справа по оси x), чтобы сыграть раунд. Получившийся результат приведет к изменению построенной кривой распределения соответствующей машины.

Естественно, чем больше раундов мы сыграем, тем точнее будут кривые распределения и тем точнее будет оценка истинного ожидания выигрыша. После того как мы сыграем определенное количество раундов, кривые распределения станут намного более точными (рис. 7.17).

Как и в случае с алгоритмом верхней доверительной границы, у машины с более высокими истинными ожидаемыми выигрышами будут более точные кривые распределения. Причина этого в том, что алгоритм работает так, что больше использует лучший автомат.

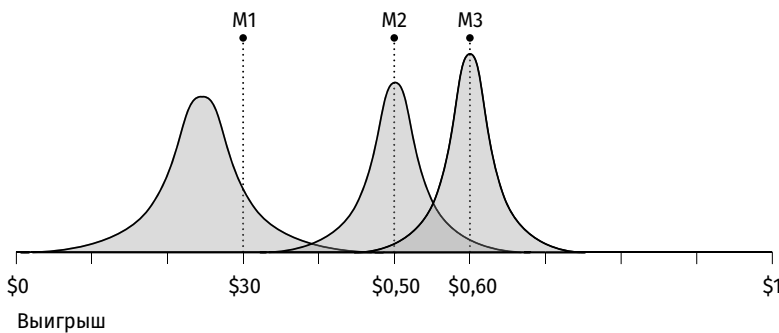


Рис. 7.17. Выборка Томпсона для уточненных диаграмм распределения вероятностей выигрыша

Тем не менее еще остается поле для исследования. Иногда даже машина M1 может выдать наилучшее мнимое истинное ожидание выигрыша. Однако это редкое явление.

Алгоритм выборки Томпсона подходит для:

- поиска наиболее эффективных каналов сбыта;
- обработки большого количества данных о клиентах, чтобы определить наиболее эффективную рекламу.

Верхняя доверительная граница vs выборка Томпсона: что предпочтительнее?

Не существует лучшего и худшего алгоритма решения задачи о «много-руком бандите». Однако и ВДГ, и выборка Томпсона имеют свои индивидуальные преимущества.

Основное отличие двух алгоритмов — способ, которым каждый из них выбирает вариант для тестирования. Верхняя доверительная граница детерминирована. Это делает алгоритм очень прямолинейным — как только сыграем один раунд, мы используем его данные для изменения границ одного из наших вариантов. Затем продолжаем тестировать вариант, который кажется оптимальным, пока данные не приведут к снижению его верхней границы до точки ниже другого варианта.

А вот выборка Томпсона вероятностна. Случайно выбирая мнимые ожидаемые выигрыши из распределений в каждом раунде, она прогнозирует, где может быть фактический результат для каждого из «бандитов», и выбирает оптимальный автомат в соответствии с этим предположением. Таким образом, в каждом раунде мы могли бы сыграть на любом автомате; нет способа сделать правильный выбор, пока не будут отображены мнимые ожидаемые выигрыши.

У обоих алгоритмов есть свои особенности. Верхняя доверительная граница обновляется после каждого раунда. Это означает, что данные, которые вы получаете, должны быть включены в значения границы прежде, чем вы сможете перейти к следующему раунду. Если вы не настроите свои значения на основе собранных данных, алгоритм не станет работать: никакие новые данные не означают, что следующий раунд будет

идентичен предыдущему, и поэтому алгоритм ничему не научится. Однако при выборке Томпсона алгоритм может учитывать отложенную обратную связь — даже если вы нечасто обновляете кривые распределения за счет собранных данных. Из-за своего вероятностного характера алгоритм будет продолжать генерировать гипотетические ожидаемые выигрыши, проверять машины и учиться на раундах.

Это одна из особенностей, составляющих преимущество выборки Томпсона: нам не нужно обновлять алгоритм данными каждый раз, когда мы играем раунд. Это не важно в случае игровых автоматов, потому что если мы будем манипулировать рычагами, то сыграем немного раундов и соберем небольшой объем данных. Если, однако, мы используем этот алгоритм для сайтов и рекламных объявлений — где могли бы получать тысячи кликов в день и где работать со всеми этими данными было бы затруднительно, — отсутствие необходимости обновлять алгоритм после каждого раунда является плюсом. Таким образом, выборка Томпсона позволяет обрабатывать данные партиями, а не каждый раз, когда появляются новые данные. В случае с сайтами это означает, что мы можем подождать до тех пор, пока не наберем определенное количество кликов, и только потом обновить алгоритм с этой информацией. Это называется пакетной выборкой.

Конец теста А/В

В начале этой главы мы коснулись часто критикуемого теста А/В. Теперь, когда вы знаете о верхних доверительных границах и алгоритмах выборки Томпсона, я надеюсь, вы видите, что для проведения подобных исследований доступны гораздо более мощные (и простые) инструменты.

Тесты А/В фокусируются на исследовании и имеют минимальную значимость для использования. Тестирование всех вариантов проводится одинаково, поэтому использование оптимального варианта возможно лишь в ограниченной степени. Однако ВДГ и выборка Томпсона ориентированы на использование. Они разработаны так, что производят только самые необходимые исследования и больше всего будут заняты лишь одним оптимальным вариантом.

Я считаю, что для небольших специальных исследований вы все равно можете применить тестирование А/В, если это вам удобно. Но в случае серьезных бизнес-проектов, особенно при значительном финансировании, ВДГ и выборка Томпсона работают гораздо более эффективно.

Будущее анализа данных

Мы хорошо потрудились, чтобы достичь конца этого раздела; самая сложная часть позади. Если вы прочитали эти две части про анализ данных, то теперь хорошо подготовлены для того, чтобы *интуитивно* понимать, какой тип алгоритма — от классификации и кластеризации до обучения с подкреплением — нужно использовать, чтобы наилучшим образом решить поставленную задачу.

Теперь сделаем передышку, прежде чем переходить к заключительной части; вы ее заслужили. Подготовка и анализ данных на сегодняшний день являются, безусловно, самыми техническими стадиями анализа и обработки данных, но если вы завершите их, то подойдете к двум последним этапам процесса во всеоружии.

ЧАСТЬ ТРЕТЬЯ

«Как я могу это показать?»

Представление данных

Сейчас мы на финишной прямой и, как я упоминал в конце предыдущей главы, вы должны чувствовать себя в состоянии легко управлять данными, если потратили достаточно времени на их надлежащую подготовку и анализ. Тем не менее не стоит расслабляться: конечные этапы анализа и обработки данных необязательно пройдут гладко. Здесь нужны совсем не те инструменты, что использовались на первых стадиях. Теперь от вас потребуются творческое мышление.

На двух заключительных этапах процесса анализа и обработки данных следует ориентироваться на других участников проекта, чтобы четко и доходчиво донести до них информацию, полученную в ходе нашей работы. В этой части я расскажу о многих способах, позволяющих аналитикам данных как можно лучше представить результаты проекта; я дам советы по визуализации и презентации, а также несколько рекомендаций относительно того, как создать вашу собственную нишу в этой области.

Хорошее впечатление

Смысл визуализации данных не только в том, чтобы делать красивые картинки. Существуют надежные методы, с помощью которых вы

можете добиться того, чтобы визуальные эффекты подкрепляли ваши идеи и впечатлили значимых людей. В главе 8 я расскажу о методах убеждения представителей всех заинтересованных сторон.

Поскольку визуализация частично зависит от аналитики, в эту главу я также включил визуальный анализ. Не только из-за сходства названий, но и потому, что визуальный анализ часто считается промежуточным звеном между анализом данных и их визуализацией. Этот тип анализа может быть использован в обоих случаях: либо перед третьим этапом (анализ данных) — с целью выработать наилучший путь построения алгоритмов, либо перед четвертым этапом (визуализация) — с целью сделать результаты более весомыми.

Вы еще не закончили!

Некоторые ошибаются, думая, что получение результатов означает завершение работы. Эти люди не понимают, как важно научиться представлять свои данные. Вот почему процесс анализа и обработки данных состоит из пяти, а не из четырех этапов. И это тема главы 9.

Все мы разные. Вам может быть интереснее заниматься техническими аспектами (алгоритмы и программное обеспечение) или осмысливать то новое, что вы обнаружили. В любом случае ошибочно пренебрегать процессом презентации. Это может прозвучать резко, но, если вы не научились хорошо преподносить результаты своей деятельности, значит, вы не отвечаете всем критериям, необходимым для того, чтобы быть топ-аналитиком данных. Вы всего лишь квалифицированный жонглер числами, однако таких людей немало вокруг и относительно немногие из них могут хорошо презентовать полученные результаты. В главе 9 я покажу вам свои лучшие уловки, которым я научился, делая презентации, и объясню, как вы можете применить их в вашей деятельности.

Выбор карьеры

Эта заключительная глава вернется к предложенным в первой части концепциям. В первом разделе мы обсудим, какие области доступны

читателям, желающим построить карьеру в области науки о данных, и какие возможности предлагает эта сфера. Читателям, которые должны подготовиться к собеседованию, я расскажу, о чем их могут спросить и как произвести впечатление на работодателей.

Наконец, мы также рассмотрим, как разработать убедительную стратегию включения науки о данных в деятельность компании. Вы можете работать в организации, недостаточно использующей данные, или быть предпринимателем, заинтересованным в разработке плана реализации собранных данных. В этом разделе раскрывается, как сделать науку о данных необходимостью, а не роскошью, тем самым обеспечив себе долгую и успешную карьеру.

Говоря о визуальных инструментах в науке о данных, можно выделить два направления: визуальный анализ и визуализацию. Разница между ними важна. Воспринимайте *визуальную аналитику* как дополнительный инструмент для этапов 1–3 в процессе анализа и обработки данных (выявление вопроса, подготовка и анализ данных). *Визуализация* — это то, что лежит в основе этапа 4, визуализации данных.

В этой главе мы научимся понимать различия между этими двумя направлениями и узнаем, как они могут улучшить наши проекты. По сути, визуальная аналитика — это изнанка: она предназначена исключительно для нас и нашей команды, чтобы мы могли тщательно изучить результаты, в то время как визуализация — лицевая сторона, мощный способ представить наши выводы всем участникам проекта.

Ниже мы расскажем о наиболее полезных инструментах визуализации и визуальной аналитики на рынке, о том, почему аналитикам данных необходимо изучать эти технологии, а также об основных стратегиях, обеспечивающих использование самых действенных визуальных эффектов в презентациях.

Что такое визуальный анализ данных?

Визуальный анализ технически устраняет разрыв между анализом данных и визуализацией, поскольку он заимствует принципы из обеих областей и может быть проведен до, во время или после обработки данных. Проще говоря, визуальный анализ помогает нам «видеть» данные — выявлять тенденции и аномалии в наших записях как бы с высоты птичьего полета. С помощью визуального анализа мы помещаем подготовленные данные в интерактивные объекты, диаграммы

и графики, позволяющие уяснить, куда ведут тенденции и аномалии. Огромное преимущество состоит в том, что мы можем применить визуальный анализ в любой момент процесса и тем самым тщательно изучить, как данные отвечают на наши вопросы.

Хотя визуальный анализ, возможно, в конечном итоге упрощает проблему, он также может помочь определить отправную точку или обобщить результаты более подробного исследования. Работая с большими объемами данных, человеку трудно разобраться в них только с помощью цифр и таблиц. Как только мы придаем этим данным некую форму и добавляем к ним цвет, направление и движение, они становятся намного понятнее.

Прежде чем перейти непосредственно к визуализации, давайте немного поговорим об аналитике.

Визуальный анализ: до или после анализа данных?

Мне нравится применять принцип Парето, чтобы проиллюстрировать, как визуальная аналитика может способствовать процессу анализа и обработки данных. Согласно этому принципу 80% последствий являются результатом 20% причин — концепция, часто применяемая в бизнесе для описания того, как, например, 20% клиентов могут обеспечить 80% продаж. Визуализация наших данных должна прояснить, где находятся наиболее релевантные записи данных, и, если время поджимает, легко показать, в каком месте массива данных мы должны сосредоточить свою энергию.

Этот подход эффективно работает до начала проекта, когда мы можем выделить корреляции на ранней стадии процесса. Визуальный анализ способен дополнительно помочь пересмотреть то, что мы уже обнаружили, и расширить контекст полученных результатов уже после завершения анализа. Это исключительно полезно для тех, кто имеет дело с одними и теми же массивами данных на долгосрочной основе: часто проводя визуальный анализ записей, мы можем быстро идентифицировать изменения и возникающие тенденции, не обращаясь к подробным алгоритмам.

Остается вопрос: визуализировать до или после анализа? Я рекомендую новичкам в науке о данных серьезно рассмотреть возможность

визуализации данных как до, так и после проведения надлежащего анализа. В конце концов, когда мы полностью погружены в проект, распознать тенденции может быть сложно. И именно это Меган Патни, менеджер по развитию Mike's Hard Lemonade Co (американского производителя напитка), делает, чтобы отслеживать продажи продукции своей компании.

Кейс: Mike's Hard Lemonade Co

Меган еженедельно визуализирует большие массивы данных своей компании из более чем двух миллионов единиц, чтобы показать, каково положение с запасами и продажами продукции по всей территории Соединенных Штатов и в каждом штате. Визуализации в программе Tableau (см. раздел об использовании Tableau) помогают Меган не только сообщать о производительности компании заинтересованным сторонам, но и указывать отстающий штат, тем самым подталкивая команду к принятию мер по улучшению ситуации.

Регулярно визуализируя данные компании, Меган также выяснила, где было бы полезно собирать дополнительную информацию. Часто программное обеспечение визуализации имеет опции фильтрации данных и выделения ключевых тенденций.

Фильтрация позволяет аналитикам данных просматривать отдельные категории и то, как они работают по отношению к другим переменным в данных, переходя от общего представления к очень конкретному — и все это одним нажатием кнопки.

Создавая еженедельно визуальные отчеты о данных, Меган однажды обнаружила пробел в информации. Она уже собрала данные из магазинов Walmart, но увидела, что не геокодировала принадлежавшие Walmart магазины Sam's Club (сеть розничной торговли только для членов клуба), где Mike's Hard Lemonade Co также продает свои продукты. Меган потратила время на сбор данных и ввела их в массив данных с Walmart. Теперь она смогла узнать, насколько близко магазины расположены друг к другу:

«Во-первых, я геокодировала все магазины, чтобы выяснить, у каких есть общая парковка. Затем взяла данные о продажах, чтобы понять наши результаты по этим магазинам. Выяснилось, что демопродажа в магазинах Sam's Club улучшила наши продажи в расположенном рядом Walmart. Так мы

смогли сделать вывод, что продажи в магазине Sam's Club, находящемся около Walmart, могут принести пользу, — и мы узнали это благодаря визуализации данных».

(SuperDataScience, 2016)

С помощью программного обеспечения визуализации Меган и ее команда выявили еще одну интересную особенность. Меган знала, что Sam's Club и Walmart часто находятся рядом друг с другом, иногда даже имея общую парковку. С помощью визуализации Меган обнаружила закономерности, связанные с продажами в магазинах, расположенных близко друг к другу. Это позволило ей синхронизировать продажи продуктов в франшизах Walmart и Sam's Club, деливших парковочные места.

Использование Tableau

Программное обеспечение для бизнес-аналитики Tableau позволяет брать данные непосредственно из базы, чтобы мгновенно создавать красивые информативные визуализации и представлять полученные результаты в интернете. Полная версия Tableau не бесплатна (на момент написания статьи за персональную лицензию надо было заплатить \$999), но, скорее всего, вам придется иметь с ней дело, поскольку она широко применяется во всем мире в области анализа данных.

Существует также бесплатная пробная версия, Tableau Public*. В этой версии программы некоторые функции ограничены (например, возможность сохранения файлов непосредственно на вашем компьютере), что делает ее неидеальной для корпоративного использования. Тем не менее Tableau Public для начала вполне подходит — просто убедитесь, что вы не работаете с какими-либо конфиденциальными данными компании.

Если вы хотите потренироваться работать с Tableau, помните, что можете использовать любой из публичных массивов данных, перечисленных во введении к второй части этой книги. Когда вы загружаете

* Еще одна совершенно бесплатная альтернатива Tableau — Microsoft Power BI, инструмент, который предлагает аналогичные функции; но лично я предпочитаю первую программу.

массив данных в программу, Tableau дает вам возможность перетаскивать различные визуальные инструменты (например, географические карты, круговые диаграммы, точечные диаграммы) и тем самым сегментировать информацию. Графики и диаграммы могут быть дополнительно украшены цветом, который Tableau определит интуитивно, а также вам будет доступен широкий диапазон вариантов форматирования для выделения сегментов с помощью ярлыков. Лучший способ учиться — это практика, и благодаря Tableau вы сможете организовать собственное обучение: выделить некоторое время для себя в выходные дни, загрузить массив данных в программу и заняться визуализацией.

Использование визуального анализа для начала анализа

Давайте представим, что вас вызвали в качестве аналитика данных в общенациональный немецкий банк, который уже определил свою проблему: нужно пересмотреть и в конечном итоге изменить то, как банк рекламирует и продает страховые продукты своим клиентам, и сегментировать их в соответствии со сферами его интересов.

Вы могли бы сначала рассмотреть, есть ли какие-либо очевидные категории, чтобы сгруппировать людей на основе подготовленных банком данных.

Как мы видели в главах 6 и 7, для решения этой проблемы можно использовать множество различных подходов. С помощью визуального анализа, чтобы добавить цвет (в буквальном и переносном смысле) к нашим выводам, мы также можем получить любые демографические данные, которые банк собрал о своих клиентах, — возраст, почтовый индекс, арендная плата, лояльность, средние расходы, размер семьи — и изучить логические выводы и соотношения между этими демографическими данными. Добавляя данные в Tableau, Power BI или любое другое программное обеспечение визуализации, мы можем получить некоторые важные сведения, которые в противном случае было бы непросто увидеть. Прогоняя данные через инструмент визуализации, мы выделяем демографию и получаем представление о том, как данные функционируют в условиях сочетания различных демографических характеристик.

Моя первая попытка, учитывая специфику банковского дела, состояла бы в разделении клиентов по их возрасту. Я могу сделать логичное предположение, что деление по возрасту часто используется в банковской отрасли, потому что банковские кредиты (для студентов, для покупателей первого дома, для молодых семей) зависят в первую очередь от возраста клиентов. Инструмент визуализации данных может показать нам, как возраст влияет на другую демографическую информацию, предоставленную банком.

Использование дашбордов/панелей индикаторов

Программы визуальной аналитики помогут изучить взаимосвязь переменных возраста с другими переменными. Если мы используем Tableau или Power BI, эту взаимосвязь можно даже визуализировать на интерактивном дашборде, который сгруппирует наших клиентов в сегменты на основе поведения и демографии.

Предположим, что наши визуальные данные говорят, что живущим в Саксонии клиентам банка в основном от 45 до 55 лет и большинство из них имеют низкие доходы. Однако клиентам банка, живущим в соседней федеральной земле Бавария, как правило, от 25 до 35 лет, и они хорошо зарабатывают. Это означает, что визуальные эффекты привлекли внимание к двум соответствующим демографическим данным в дополнение к возрасту: местоположению и доходу. Теперь в Баварии мы могли бы сделать акцент на выгодных для молодых клиентов льготных кредитах — на приобретение первого дома, на свадьбу и семью. В Саксонии мы могли бы ориентироваться на защиту пожилых клиентов — страхование жизни, имущественное планирование и инвестиции с низким уровнем риска.

Мы можем, конечно, продолжать и дальше сегментировать наши данные, как считаем нужным, но выявленные подкатегории уже добавили ценность нашему проекту благодаря визуальному анализу данных. Это оказалось гораздо проще, чем использовать аналитические инструменты, которые мы рассмотрели в предыдущих главах: потребовалось только программное обеспечение для преобразования структурированной информации в визуально привлекательный формат.

В итоге банк может не только рассылать брошюры, соответствующие возрасту потребителя, но и советовать своим маркетологам адаптировать материалы для дополнительных подкатегорий: с низким доходом, высоким доходом и т. д.

Здесь мы увидели пример того, как визуальная аналитика может помочь на третьем этапе работы над нашим проектом, а также как можно продолжать использовать ее для отслеживания новых тенденций, расхождений и пробелов в данных. На следующем этапе можно будет перейти к визуализации данных, чтобы их могли оценить те, кто наиболее важен для нашего проекта: участники и другие заинтересованные стороны.

Суть визуализации данных

К настоящему времени вам уже должно быть ясно, насколько полезна визуализация данных для их исследования. Итак, просто представьте, какую существенную роль могут сыграть визуальные эффекты при представлении нашей информации участникам проекта.

Визуализация данных — это процесс создания наглядных средств, помогающих людям видеть и понимать информацию. Визуализируя данные, мы представляем их в контексте. Ведь без контекста данные бессмысленны. Нам нужно задать данным правильно сформулированный вопрос. Мы должны подготовить наши данные в стандартизированном формате, а затем проанализировать их. Все эти шаги способствуют добавлению контекста, необходимого для понимания предоставленной информации. Но в то время как контекст помогает распознавать тенденции и результаты данных, мы должны продолжать расширять его, чтобы заинтересованные стороны оценили всю нашу тяжелую работу. Ведь мы не можем познакомить участников проекта только с результатами анализа — маловероятно, что они сведущи в науке о данных. Кровь, пот и слезы, которые вы вложили в подготовку и анализ данных, скорее всего, для них ничего не значат. Не принимайте это на свой счет.

Все, что нам нужно сделать сейчас, — это обеспечить перевод наших результатов в то, что может быть легко понято заинтересованными

сторонами. Здесь мы можем чему-то научиться у наших друзей из мира бизнеса. BI* в первую очередь связана с составлением и представлением отчетов, которые помогают улучшить и изменить работу бизнеса. Проведение бизнес-анализа может показаться несложным, учитывая нашу предыдущую главу об анализе данных, но я не знаю ни одного руководителя, который предпочел бы прочитать 50 страниц, вместо того чтобы изучить визуально представленные проблему и пути ее решения. При качественной визуализации дашборды BI привлекут и убедят аудиторию внести предложенные вами изменения.

Здесь мы также должны позаботиться о том, чтобы понять, как работают визуальные эффекты — у них есть свой собственный язык, которым вам необходимо овладеть, иначе ваши презентации могут скорее запутать, чем убедить участников.

Важность визуальности

Вот мысль, которая, как я знаю, расстроит многих аналитиков данных: в некотором смысле ваш проект не касается ни данных, ни визуальных элементов. В конечном счете речь идет о *людях*. Нас привлекли к решению бизнес-вопроса, который повлияет на *заинтересованных лиц*, будь то клиенты или руководители. И если наша информация не будет представлена таким образом, чтобы ее поняли люди, которые могут санкционировать изменения, рекомендованные нашими данными, — тогда все наши усилия окажутся потрачены впустую.

Проще говоря, результаты в электронной таблице видны хуже, чем на диаграмме. Но это не значит, что мы должны просто держать наготове набор красок; даже использование такого распространенного средства, как цвет заливки таблиц в Excel, может не сработать.

* BI (Business intelligence) — IT-технологии для сбора, хранения и анализа данных. На основе информации, собранной и проанализированной с помощью BI, можно принимать эффективные решения для управления бизнес-процессами. — *Прим. науч. ред.*

Я бы сказал, что это одна из причин того, почему многие реализованные проекты в области науки о данных, к сожалению, никуда *не приводят*: слишком много практиков считают, что данные будут говорить сами за себя. К сожалению, это не так. Нам надо научиться убеждать нашу аудиторию, а визуализация — это канал, по которому непременно должны пройти данные, иначе наши озарения рискуют навеки застрять в чистилище компании.

Всегда ли нужна визуализация?

Некоторые спрашивают меня, должны ли аналитики данных всегда использовать визуальные элементы в своих проектах. Однозначный ответ — нет, не *всегда*. Нам нужно визуализировать данные, только когда идеи *не могут быть эффективно представлены без визуальных эффектов*. Но по указанным выше причинам визуальные элементы почти во всех случаях будут полезны для проектов, требующих изучения массивов данных. Вам редко придется работать над тем, чтобы дать только один прямой ответ относительно лишь одного аспекта данных, — в таких случаях визуализация не является строго необходимой.

Говорить на визуальном языке

Визуальные элементы передают информацию способом, недоступным тексту. Восприятие наших *написанных* аргументов может быть ограничено многими факторами, в частности скоростью чтения, уровнем понимания вопроса и наших рекомендаций. А единичное изображение может проиллюстрировать особенности, выявить отклонения и определить отдельные группы данных, эффективно минуя сотни строк массива данных для информирования нашей аудитории — и все это будет понятно буквально с одного взгляда.

Еще одно преимущество использования изображений в этой области связано с тем, что мы живем в очень визуальной культуре.

Мы привыкли к картинкам в социальных сетях, к интернет-мемам — и наш аппетит к общению посредством изображений увеличился. Инфографику и другие визуальные средства представления информации освоили учреждения в мире бизнеса, которые хотят наглядно показать своей аудитории, что они делают. Таким образом, изображения больше не являются декоративным дополнением к презентации в PowerPoint; люди теперь *ожидают* визуальную информацию. Поэтому, используя изображения для представления данных, вы просто реализуете эти ожидания*.

В результате этих перемен мы овладели визуальным языком. Впитывая культуру современного общества, мы укоренили в себе знания о том, как могут быть интерпретированы различные цвета и формы. Вы удивитесь, увидев, как много уже знаете из того, о чем здесь пойдет речь дальше. И хотя следует учитывать, что визуальный язык не всегда является международным, его часто можно использовать для пользы дела.

Самообслуживающаяся аналитика

К счастью, разработчики отреагировали на растущие потребности в мощных визуальных средствах, придумав программное обеспечение и приложения для создания интерактивных визуализаций. С помощью таких программ, как Tableau и PowerBI, мы можем создавать визуальные эффекты, которые показывают широкую картину и — благодаря тщательному использованию фильтров — позволяют увидеть детали (для этого достаточно нажать кнопку, выбирая интересующую категорию). Будучи динамичными, эти программы уменьшают нагрузку на аналитиков данных, так как конечные пользователи могут управлять визуальными элементами так, чтобы представить данные именно в том виде, в каком нужно.

* Вдохновляющие идеи для визуализации проектов смотрите на сайте журналиста, работающего с данными, Дэвида Маккэндлесса, informationisbeautiful.net, что выводит визуальное повествование на новый уровень.

Есть также приложения, доступные для самых популярных программ анализа данных, — R и Python. С их помощью вы можете создавать визуальные средства прямо в ходе проведения анализа, не выходя из программы. Seaborn — бесплатный пакет визуализаторов данных для Python; он дает возможность пользователям рисовать статистическую графику по своим данным. Plotly позволит вам создавать интерактивные визуализации для самообслуживающегося анализа. В программе R то же самое может быть достигнуто с помощью ggplot2 и Shiny.

С развитием науки о данных и распространением основанного на данных подхода к бизнесу тенденция к самообслуживающейся аналитике усилилась. Побуждение других к самостоятельному поиску ответов на вопросы полезно всем, и именно это движет мной в профессиональной деятельности. Я советую это любому специалисту по данным.

Построение привлекательных визуальных элементов

Хорошие визуальные элементы информируют аудиторию. Они помогут людям увидеть наши данные. Хотя я считаю, что визуальный язык понятен, особенно в компьютерный век, есть несколько приемов, которым вы можете научиться, чтобы более эффективно воздействовать на свою аудиторию с помощью визуальных эффектов.

1. Вернитесь к началу

Вернитесь к первому шагу процесса анализа и обработки данных: определите вопросы и имейте их в виду, когда начнете отбирать визуальные элементы. Что хочет знать ваша аудитория? Представляете ли вы визуальные материалы нескольким вовлеченным сторонам и если да, то будут ли у них разные интересы? В каких точках эти интересы пересекаются? Что участники смогут узнать из ваших данных?

Ответы на эти вопросы должны быть вам уже известны в случае, если вы провели анализ данных. У вас есть несколько результатов, существенных для причастных сторон? Расположите полученные результаты в порядке важности и убедитесь, что самый значимый из них отображен как наиболее приоритетный (занимает большую часть страницы), и/или поместите его в верхний левый угол, чтобы аудитория увидела его в первую очередь.

2. Ограничьте объем текста

Я предпочитаю использовать как можно меньше текста, позволяя визуальным элементам выполнять большую часть работы, так что мне остается только добавить «красоты» в ходе презентации.

Ваша визуализация может служить двум целям — просто нужно определить, которой из двух. Если вы будете присутствовать на презентации и сможете использовать визуальные элементы для убеждения аудитории, то чем меньше текста, тем лучше. Вам все равно понадобятся метки и заголовки для диаграмм, но цель должна состоять в том, чтобы удержать внимание аудитории. Визуальные эффекты нужны, чтобы помочь *вам*, а не кому-то другому.

Если же вы отправляете отчет или дашборд самообслуживающейся аналитики, вам может потребоваться дополнительный текст, чтобы для удобства пользования читателей объяснить, в чем суть. Эффективный способ проверить, хорошо ли составлено сообщение, — попросить членов команды, чтобы они просмотрели ваш результат и оценили, насколько он понятен.

3. Добейтесь четкости

Если у вас представлено несколько объектов, убедитесь, что между ними есть четкие границы. Вы ведь не хотите, чтобы ваша аудитория увидела связь там, где ее нет, иначе получится, что вы непреднамеренно манипулируете данными. Если визуальные элементы должны сопровождаться легендой (объяснением, добавленным к диаграмме для описания категорий), убедитесь, что она помещена в отдельное поле.

4. Направляйте аудиторию

Все визуальные элементы должны отвечать на конкретный бизнес-вопрос, поэтому часто излишним оказывается пошаговое объяснение того, как ваши выводы связаны друг с другом. Подробное разъяснение содержания визуальных материалов поможет следовать от одного визуального элемента к другому,

К счастью, мы предсказуемые существа, поэтому простое размещение информации сверху вниз и слева направо для каждого из объектов гарантирует, что ваша аудитория прочитает визуальные эффекты в правильном порядке. Хотя вы, возможно, удивляетесь примитивности этого подхода, я уверен, что вы можете вспомнить несколько примеров, когда созданный с самыми благими намерениями текст или плакат превратился в интернет-мем.

5. Приручите своего внутреннего творческого зверя

Иногда создание визуальных эффектов может привести в офис живущего в нас дизайнера, готового изменить мир и вооруженного тепловыми картами. Однако не поддавайтесь этому желанию, оставьте свои художественные пристрастия дома.

Не думайте, что игра с размерами шрифтов поможет вашей аудитории отличить самую важную информацию от наименее существенной. На самом деле произойдет обратное: большая часть присутствующих будет стремиться прочитать текст, набранный самым мелким шрифтом, обращать внимание прежде всего на то, что вы добавили в качестве сносок. Поэтому используйте два шрифта: один для заголовков, а другой для основного текста, и, если у вас есть соблазн уменьшить размер шрифта для любой информации, спросите себя, не следует ли ее вообще выбросить.

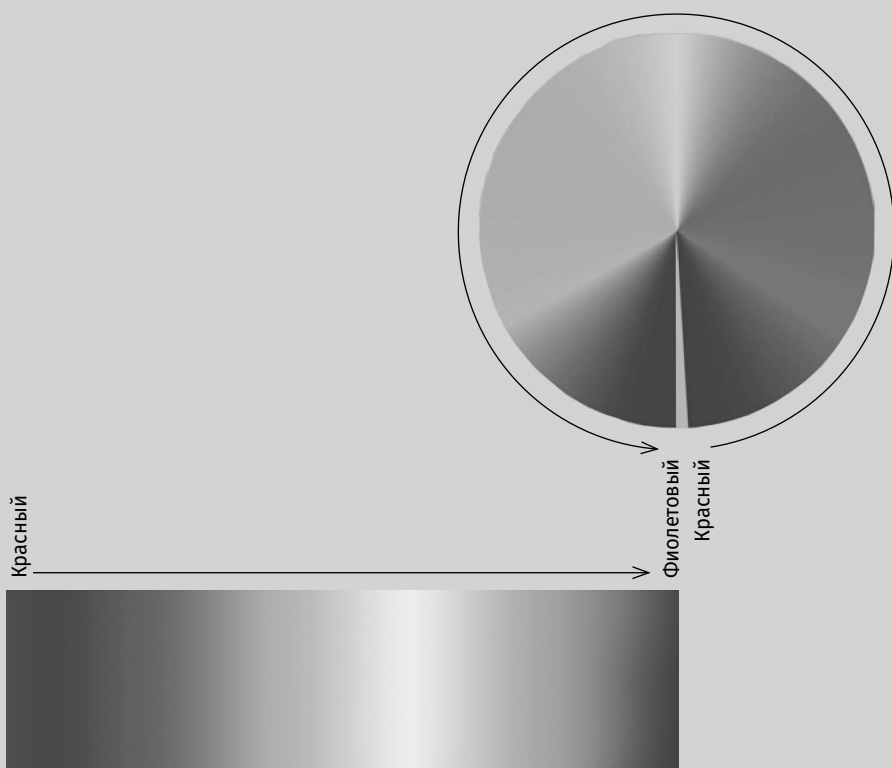
6. Будьте проще!

Если у вас есть несколько визуальных объектов, убедитесь, что они хорошо смотрятся вместе. Будьте последовательны и по возможности избегайте использования более чем одной цветовой гаммы, если

Цветные схемы

Многие люди недооценивают важность цвета. Я был одним из них, пока не изучил теорию цвета и не увидел, насколько могу улучшить свои визуализации, просто правильно комбинируя цвета.

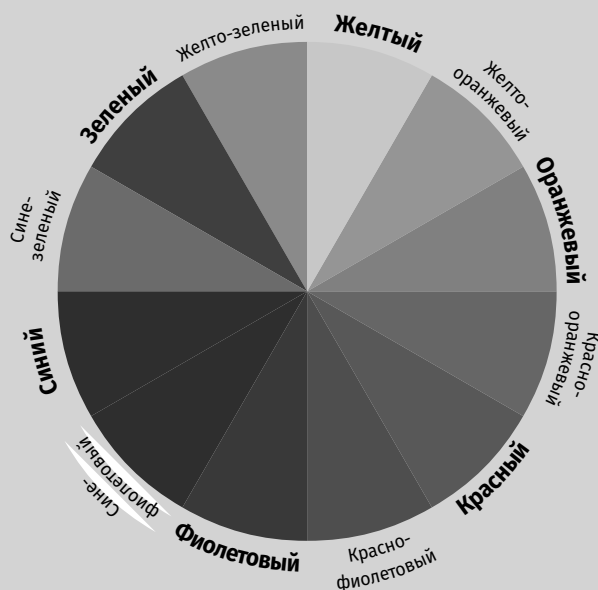
Самое приятное то, что за этим стоит наука. Вам не нужно выбирать цвета случайным образом. Некоторые цвета лучше работают сообща. Некоторые комбинации дают наиболее благоприятный эффект, когда они контрастируют, а некоторые — когда постепенно меняются с одного на другой. Теория цвета — секретное оружие в моем инструментарии науки о данных. Почему секретное? Потому что, как правило, аналитики данных не считают, что цвета важны.

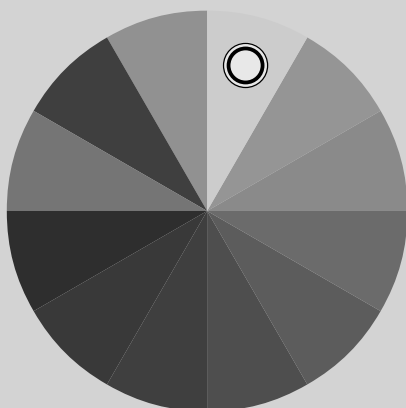


Следовать принципам использования цвета довольно просто. Лучший метод, на мой взгляд, основан на цветовом колесе. Полный спектр цветов переходит (как в радуге) от красного до фиолетового. Если эти два конца спектра соединены и образуют круг, создается традиционное цветовое колесо.

Каждому сектору круга соответствует определенный цвет или, точнее, оттенок. Оттенки незаметно переходят друг в друга. Их бесконечное множество. Для простоты я разделил колесо на шесть основных тонов: желтый, оранжевый, красный, фиолетовый, синий и зеленый. Затем я разделил шесть сегментов так, чтобы получить 12, так что, например, оттенок на полпути между желтым и зеленым получается желто-зеленым.

Далее даны шесть цветовых правил использования того же колеса. Вам нужно усвоить, какие основные элементы цветового колеса подходят друг другу, — и тогда вы сможете просто выбрать самую лучшую цветовую палитру для вашего проекта. Поскольку цветовое колесо представляет собой схематическое изображение всего цветового спектра, его черно-белый вариант может (парадоксально) помочь рационально размышлять о теории цвета, избегая субъективной реакции на сами цвета.





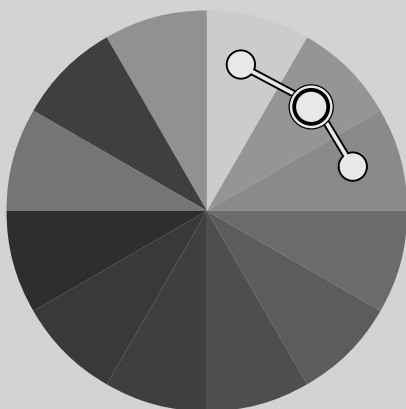
Монохромная схема

Это самая простая из схем, так как она охватывает только один цвет.

После того как вы выбрали один цвет, найдите более светлый или более темный оттенок для различения категорий. Преимущество использования монохроматической схемы заключается в том, что почти не ограничивается количество оттенков, которые вы можете использовать.

Когда использовать

Отлично подходит, если вам надо, чтобы различия в оттенках были замечены людьми с проблемами восприятия цвета.



Аналоговая схема

Выберите один цвет и два прилегающих к нему цвета. Поскольку разница между цветами здесь незначительная, эта схема может быть трудна для восприятия.

Когда использовать

Отлично подходит для тепловых карт, а также для иллюстрации постепенных изменений в данных.



Контрастная схема

Выберите два цвета, диаметрально противоположные друг другу на колесе, — желтый и фиолетовый, или зеленый и красный, или синий и оранжевый.

Когда использовать

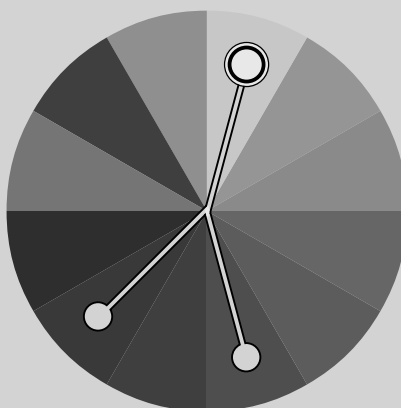
Контрастные цвета, подобные этим, удивительно эффективны при сравнении двух значений и обычно используются для выделения отдельных категорий.

Расщепленно-контрастная схема

Как следует из названия, здесь используется контрастная схема, но с одним дополнительным цветом. Вместо того чтобы брать диаметрально противоположный цвет, возьмите два оттенка по обеим сторонам от противоположного цвета.

Когда использовать

Идеально подходит, если у вас есть три категории, но вы хотите, чтобы ваша аудитория сосредоточилась на одной, в частности, для того, чтобы выделить тенденции или особенности.

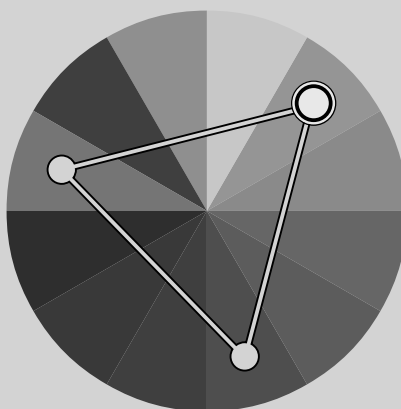


Триада

Выберите три цвета, расположенные на основной цветовой схеме в вершинах равностороннего треугольника. Например, в местах, где находятся красный, голубой и желтый цвета.

Когда использовать

Идеально подходит для изображения трех категорий равной значимости, когда нет необходимости выделить одну из них (то есть здесь расщепленно-контрастная схема не подошла бы).

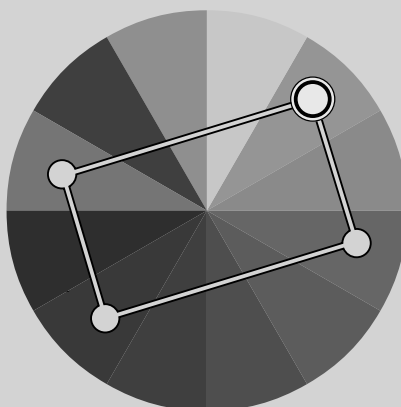


Квадратическая схема

Схожа с триадой — с той разницей, что здесь присутствуют четыре цвета, а не три.

Когда использовать

Если у вас есть несколько категорий, которые одинаково важны для ваших целей, используйте триадную схему для нечетного числа категорий и квадратическую — для четного.



только вы сознательно не берете контрастные цвета, чтобы выделить определенную часть информации. При этом старайтесь в диаграммах не окрашивать в схожие цвета несвязанные объекты, потому что ваша аудитория неизбежно воспримет их как связанные.

Более глубокие сведения о том, как можно использовать цвет, изложены ниже.

Просто помните, что чем меньше — тем лучше, и это особенно справедливо в случае презентаций. Не нагромождайте объекты и не сходите с ума от цветов. Убедитесь, что ваш текст хорошо виден на цветном фоне. Придерживайтесь выбранной цветовой гаммы и охватите белое пространство — если что-то не служит вашей аргументации, отбросьте это, как горячий картофель.

7. Помните о людях с нарушениями зрения

При использовании цвета в визуальных элементах помните, что необходимо быть внимательным к тем, у кого есть проблемы со зрением. Существуют специальные палитры цветов, разработанные с учетом таких патологий. Эти палитры уже предустановлены в некоторых программах, например в Tableau.

8. Не бойтесь сокращений

Имейте в виду, что визуализация может привести к потере информации. Часто, чтобы создать оптимальный визуальный образ и тем самым сделать свое сообщение максимально эффективным, приходится объединять или обобщать информацию. Но визуальные элементы никогда не будут *добавлять* информации (если они это делают, значит, вы манипулируете данными, что недопустимо). Например, нет необходимости включать все числа из массива данных — относительные различия на диаграммах, которые вы используете, должны представлять лишь некоторые из них. Иногда названия столбцов также могут быть удалены в зависимости от того, какую идею вы хотите донести до своих слушателей.

Сведение визуальных элементов к их сути может создать гораздо более мощную и убедительную картину; просто убедитесь, что вы знаете, что за ней стоит, на случай если вам будут задавать вопросы.

В этом и заключается задача аналитика данных — найти правильное соотношение представляемой информации и визуальных эффектов.

9. Добавьте контекст

Без контекста данные бессмысленны. Визуальные средства собирают данные, иногда из самых разных источников, и это может привести к неожиданным даже для опытных практиков открытиям, связанным с контекстом. Когда мы воспринимаем вещи визуально, то можем лучше понять взаимосвязь между сведениями, которые иначе кажутся разрозненными. Визуализация обеспечивает контекст цифрам. Большинству людей визуальные элементы позволяют отдохнуть от текста в длинных отчетах. Они могут помочь сформировать первоначальные впечатления и дать вам возможность выявить аномалии и пробелы в информации (в следующей главе я покажу, как извлечь максимум пользы из аномалий в ваших презентациях).

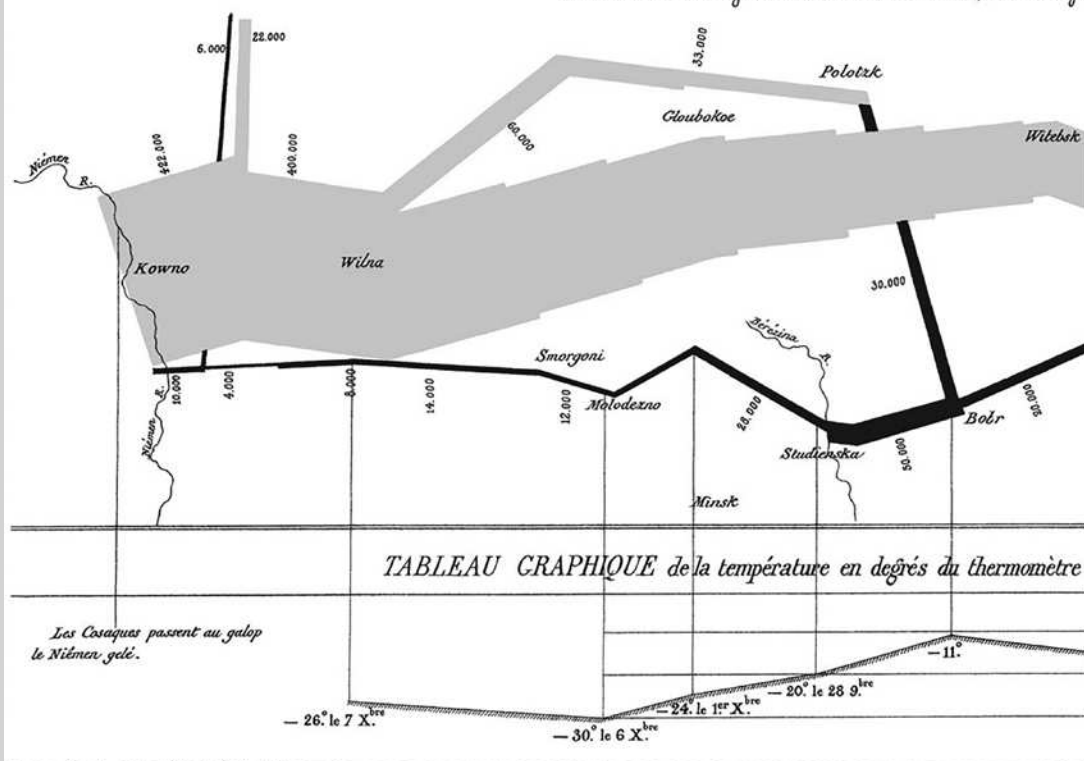
Добавление относительных значений поможет вашей аудитории понять, *почему* что-то важно. Отображенная на графике информация о тратах компании за определенный период времени может кого-то обеспокоить, но помещение этой информации в контекст (например, демонстрация размера расходов компании по сравнению с издержками конкурентов) даст возможность полностью осознать, что компания тратит слишком много из года в год.

Новости служат хорошим примером того, как информация может быть визуализирована различными способами. Предположим, что мы рассматриваем, как об индексе загрязнения атмосферы в некоей стране сообщают два новостных агентства, одно из которых работает на внутреннюю аудиторию, а другое — на внешнюю. Национальный новостной сайт мог бы сообщить, что после проведения последнего климатического саммита страна приняла серьезные меры по сокращению загрязнения и что ее индекс загрязнения атмосферы намного ниже, чем в предыдущем году. Это подошло бы для позитивного доклада. А вот заграничная новостная сеть могла с той же легкостью показать, что индекс этой страны, если принять во внимание ее относительный размер, невыгодно отличается от того же показателя в других странах,

Итак, визуализации необходим контекст.

Carte Figurative des pertes successives en hommes de l'Armée Française dans Rusie par M. MIMARD, Inspecteur Général des Ponts et C

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre des zones. Le rouge désigne les hommes qui ont été en Russie, le noir ceux qui en sont sortis. Les, dans les ouvrages de M. M. Thiers, de Ségur, de Fezensac, de Chambray et le journal inédit de Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et de Mobilow en ont rejoint vers Orscha et Witebsk, avaient toujo

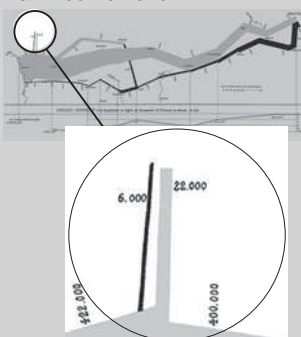


Autog. par Regnier, 9, Par. 5th Marie 5th 0th à Paris.

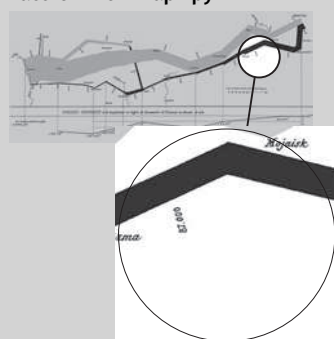
Отступление Наполеона из Москвы (Русская кампания 1812–1813 гг.)

5
слоев
информации

Количество войск



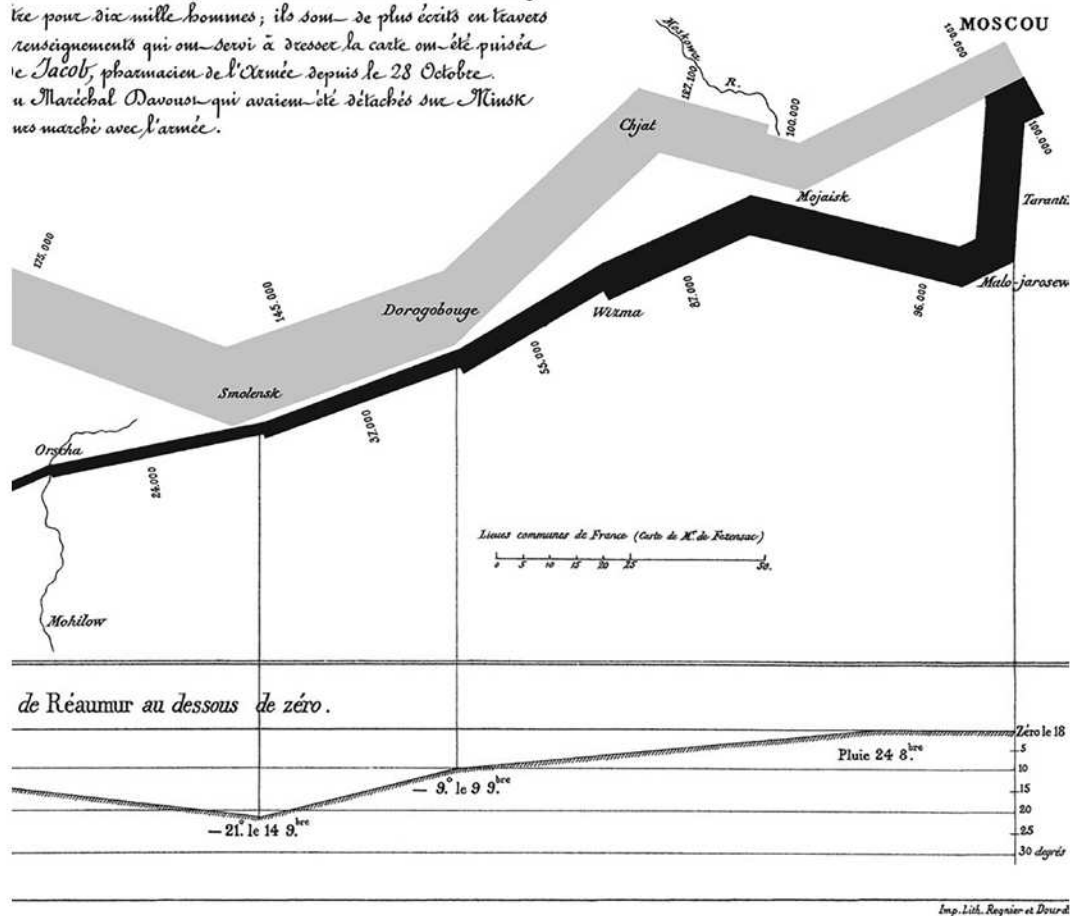
Расстояние и маршруты



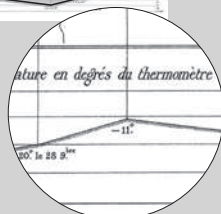
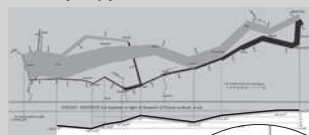
us la campagne de Russie 1812-1813.

Thausièr en retraite Paris, le 20 Novembre 1869.

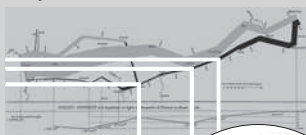
ice pour dix mille hommes; ils sont de plus écarts en travers
renseignements qui ont servi à dresser la carte ont été puisés
le Jacob, pharmacien de l'armée depuis le 28 Octobre.
u Maréchal Davout qui avaient été détachés sur Minsk
nes marchés avec l'armée.



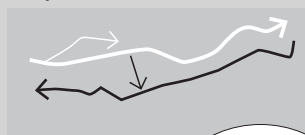
Температура



Широта и долгота



Направление движения



Слои значения

Методы визуализации информации начали формироваться в XIX в. Посмотрите на диаграмму ниже. Составленная в 1869 г. гражданским инженером Шарлем Жозефом Минаром, она показывает наступление армии Наполеона на Москву и ее отступление от Москвы в 1812–1813 гг. Здесь использован один из самых важных способов представления информации: послойный. Диаграмма Сэнки (которую мы обсудим позже) накладывается на географическую карту.

График принимает форму карты, отображающей путь великой армии из Ковно в Москву и обратно. Серая линия прослеживает ее движение к центру России, а черная — отступление. В любой заданной точке ширина линии показывает количество военных. С первого взгляда мы можем видеть разрушительные последствия кампании для французских войск. Минар говорит о том, как много солдат выступили в поход: 422 000 (число в верхнем левом углу, выше Ковно). К моменту возвращения армии ширина черной линии говорит нам о том, что число уцелевших солдат сократилось до 10 000. Мы также можем видеть, где армия разделилась, какие города она прошла, их географическое положение, какие реки они перешли вброд и какова была температура воздуха в тех местах.

Только из одного изображения мы получили огромную информацию. Минар рассказал нам историю Французской кампании в России, используя следующие слои информации:

- 1)** количество войск (задается шириной линий и числами вдоль маршрута, чтобы показать количество жертв);
- 2)** расстояние и пройденные маршруты (расстояние, указанное в легенде, маршруты, визуализированные разветвлением линий);
- 3)** температура (отображается непосредственно под картой в числовом формате);
- 4)** широта и долгота (представлены линейным графиком непосредственно под картой);
- 5)** направление движения (серая линия читается слева направо, черная — справа налево).

Эта карта показывает нам принципы изложения истории через слои информации.

10. Сохраняйте объективность

Когда мы визуализируем наши данные, немного информации неизбежно теряется. Поэтому важно не допустить фальсификации фактов при создании визуальных материалов. Это может быть трудно, особенно учитывая предпочтения и стереотипы, связанные с отдельными цветами и формами. Результат, выделенный зеленым цветом, например, может ассоциироваться с прибылью или интеллектом, тогда как в реальности ничего подобного нет.

Сказанное в основном применимо к визуальным элементам, требующим сравнительного анализа (сравнение двух или более идей, учреждений, продуктов или услуг). Как поставщики данных, мы не можем ставить одну цель выше другой из-за наших собственных предпочтений; мы должны охватить все аспекты. Участники проекта быстро увидят, что сравнительный визуальный анализ оказался однобоким, — и спросят, почему вы не включили сопоставительную информацию по всем категориям. Не скрывайте данные, которые вам не нравятся.

11. Помните о других возможных вариантах визуализации

Горькая правда состоит в том, что данные скучны для большинства людей. Существует множество вариантов визуального представления

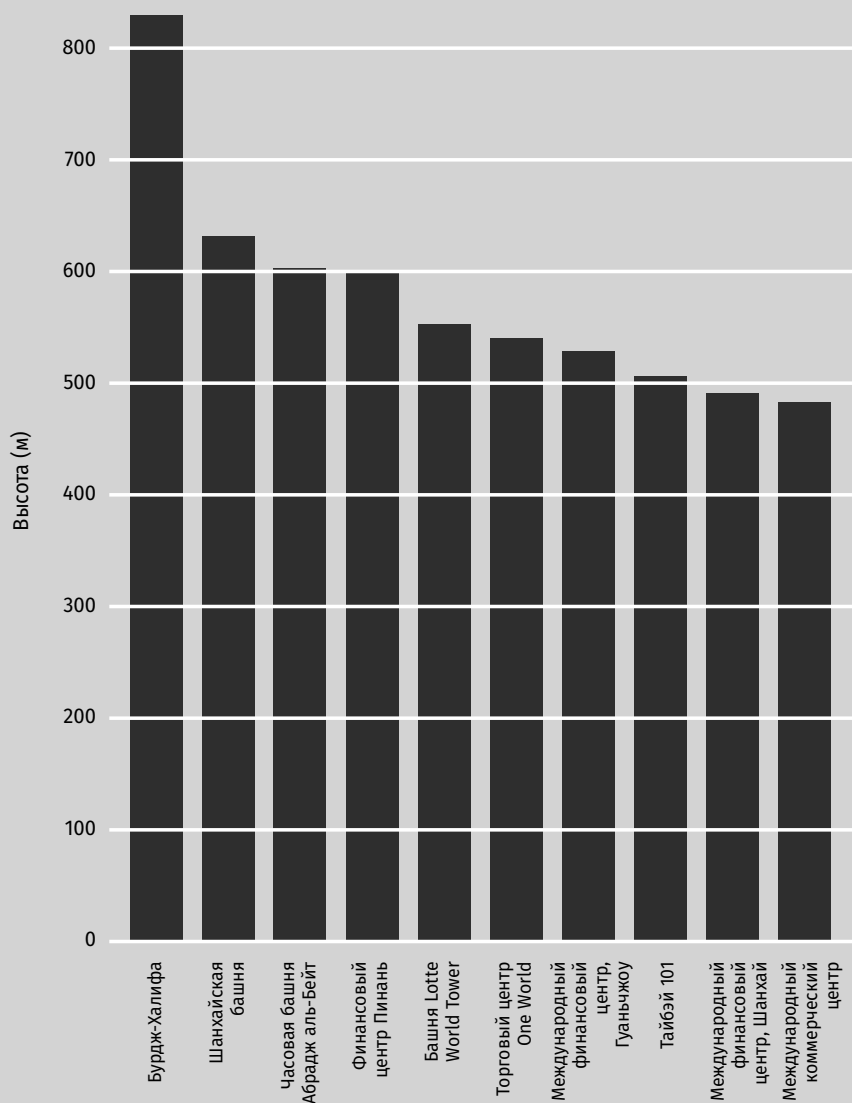
Визуализация качественных данных

Я рассказал о некоторых из моих любимых визуальных элементов, от простых до сложнейших. Вы заметите, что почти все они требуют использования количественных, а не качественных данных, чтобы быть эффективными. Визуализировать нечисловые данные может быть сложнее, хотя для этого есть много возможностей. Воспользуйтесь Sentiment Viz — генератором контента, позволяющим вводить одно или несколько ключевых слов и узнать, как пользователи Twitter относятся к ним, из твитов, отправленных за последние 20 минут*. Другая опция — облака слов, которые обсуждались в главе 3.

* Это очень интересный инструмент, и я рекомендую попробовать его. Вы можете найти его, введя «Sentiment Viz» в Google.

Гистограмма

В гистограмме данные сгруппированы в соразмерные столбцы информации, которые могут отображаться горизонтально или вертикально. Это хорошее наглядное пособие для сравнения и ранжирования количественных данных, поскольку по осям x и y значения легко различаются по длине.

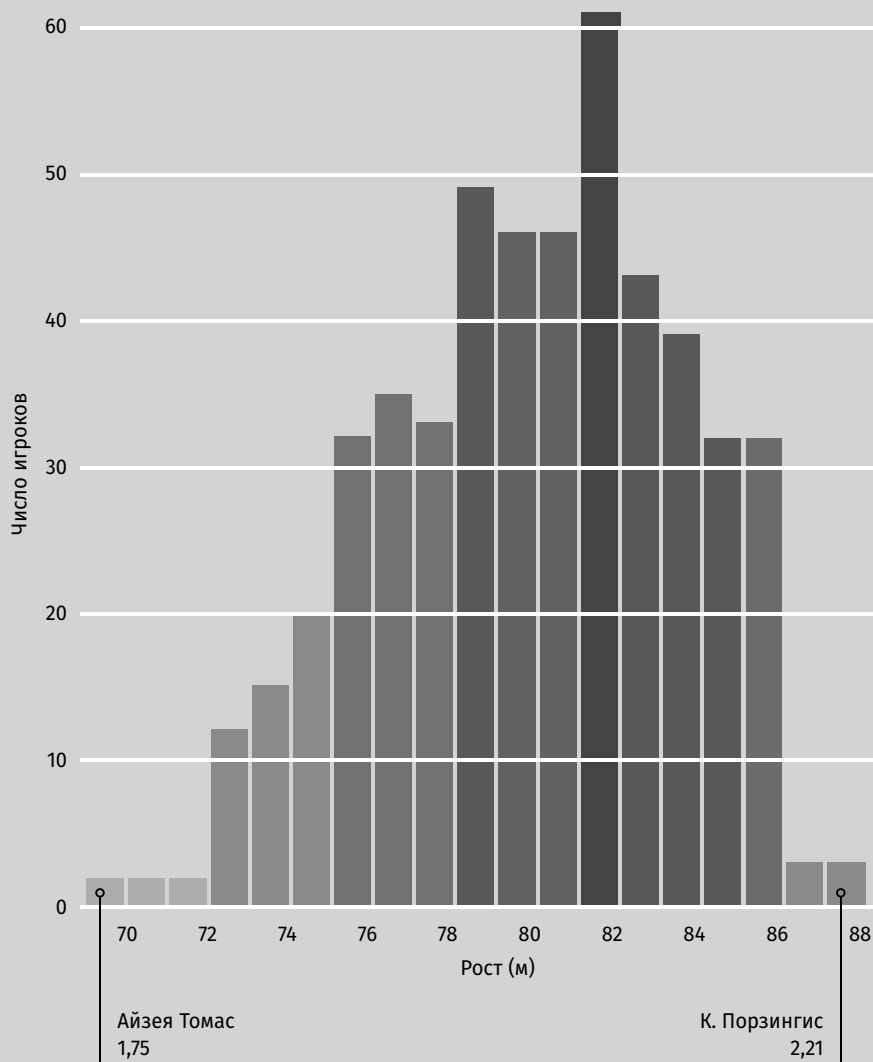


Гистограммы и распределение вероятностей

Оба варианта, показывающие, как распределяются количественные данные, подходят для выделения пиков и спадов данных, что может пригодиться в работе с вероятностями или данными переписи, такими как возраст.

На этой гистограмме мы видим распределение роста игроков НБА в сезоне 2016/17 г.

Источник данных: www.scholarshipstats.com



Линейный график

Линейный график связывает точки данных одной линией. Чаще всего он используется для отображения тенденций во временных рядах, таких как частота использования продукта в течение года.

На этом графике мы видим рост цены акций Apple за последние пять лет. Источник данных: www.data.worldbank.or

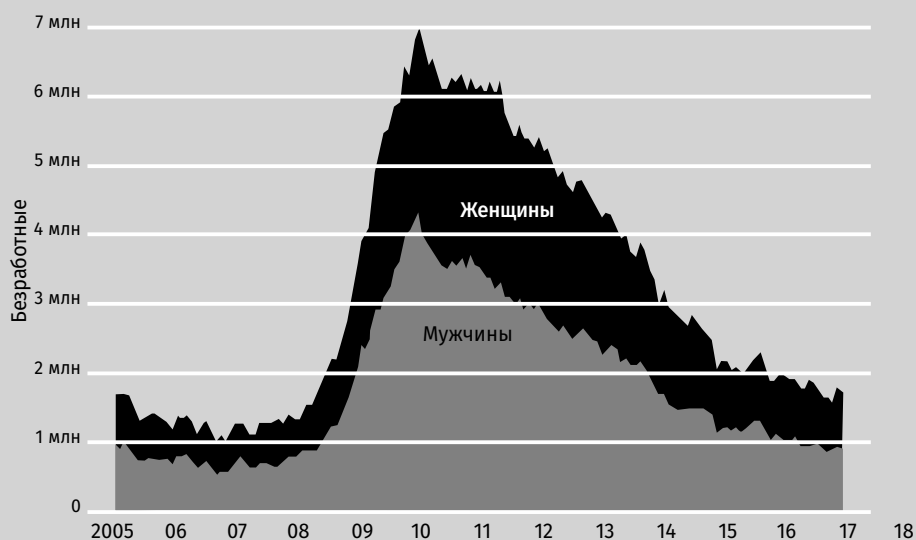


Диаграмма с областями

Диаграмма с областями — это график с зонами, выделенными цветом. Можно накладывать диаграммы областей поверх друг друга для того, чтобы категории контрастировали на рисунке. Этот тип диаграммы полезен при работе с сегментированными данными, например когда речь идет о клиентах, классифицированных по возрасту или местоположению.

На рисунке показана диаграмма, отображающая уровень долгосрочной безработицы (27 недель и более) за период с января 2005 г. по август 2017-го среди граждан США, разделенных по полу. Обратите внимание, что наложенные диаграммы также добавляются к диаграмме с областями, поэтому имеет смысл совместное рассмотрение категорий.

Данные source: www.bls.gov



Точечная диаграмма

Этот тип диаграммы помещает данные на график, основанный на двух переменных информации по горизонтальной и вертикальной осям. Расположение точек данных на диаграмме будет зависеть от их отношения к переменным.

Эта диаграмма рассеяния позволяет, например, исследовать корреляцию между рождаемостью (сколько детей рождается ежегодно на каждую 1000 существующих граждан) и процентом людей, имеющих доступ к интернету. Каждая точка представляет страну и показывает, к какой категории дохода она относится:

▲ высокий доход

* средний доход

● доход ниже среднего

▼ низкий доход

Источник данных: www.data.worldbank.org

Примечание

Это хороший пример того, что корреляция не обязательно подразумевает причинно-следственную связь. Или это не так?



Пузырьковая диаграмма

Пузырьковые диаграммы могут быть двух типов. Первый создается введением дополнительного слоя информации в точечную диаграмму — путем увеличения размера точек (и тем самым превращения их в пузырьки). Так, в нашем предыдущем примере мы могли бы присвоить размер точкам для отображения численности населения каждой страны.

Второй тип пузырьковой диаграммы гораздо менее сложный. Относительные размеры пузырьков по-прежнему кодируют данные, однако координаты отсутствуют и пузырьки расположены случайным образом.

Вот средние бюджеты фильмов в Голливуде по жанрам за 2007–2011 гг.

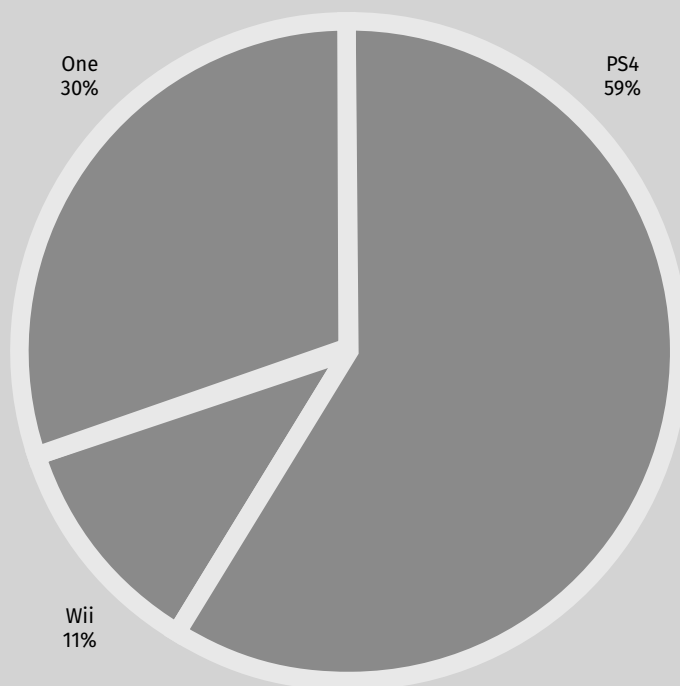


Круговая диаграмма

Многие ненавидят круговые диаграммы, потому что, в отличие от гистограмм, они не позволяют быстро увидеть разницу между категориями. При этом надо иметь в виду, что такие диаграммы могут быть очень эффективны, если: а) у вас есть сравнимые по величине данные и вы хотите показать это или б) если у вас есть сильно различающиеся по величине данные. Как правило, лучше избегать использования круговых диаграмм для отображения более трех или четырех категорий. Если элементов больше трех, такие диаграммы становятся нерепрезентативными, поскольку приходится решать, в какой последовательности расположить элементы, и в случае их неправильного порядка сопоставлять данные становится сложно.

С помощью этой круговой диаграммы можно сравнить глобальные продажи игровых консолей нового поколения за 2014–2015 гг. Платформы в этом массиве данных — PS4, One и Wii.

Источник данных: www.vgchartz.com



Плоское дерево

Плоские деревья *абсолютно* не похожи на обычные деревья. Упорядочивание данных в блоки разного размера больше похоже на пузырьковые диаграммы за исключением того, что древовидные карты немного более организованны.

На этой карте представлена первая десятка стран с самыми большими военными бюджетами (оценки 2017 г., \$млрд).

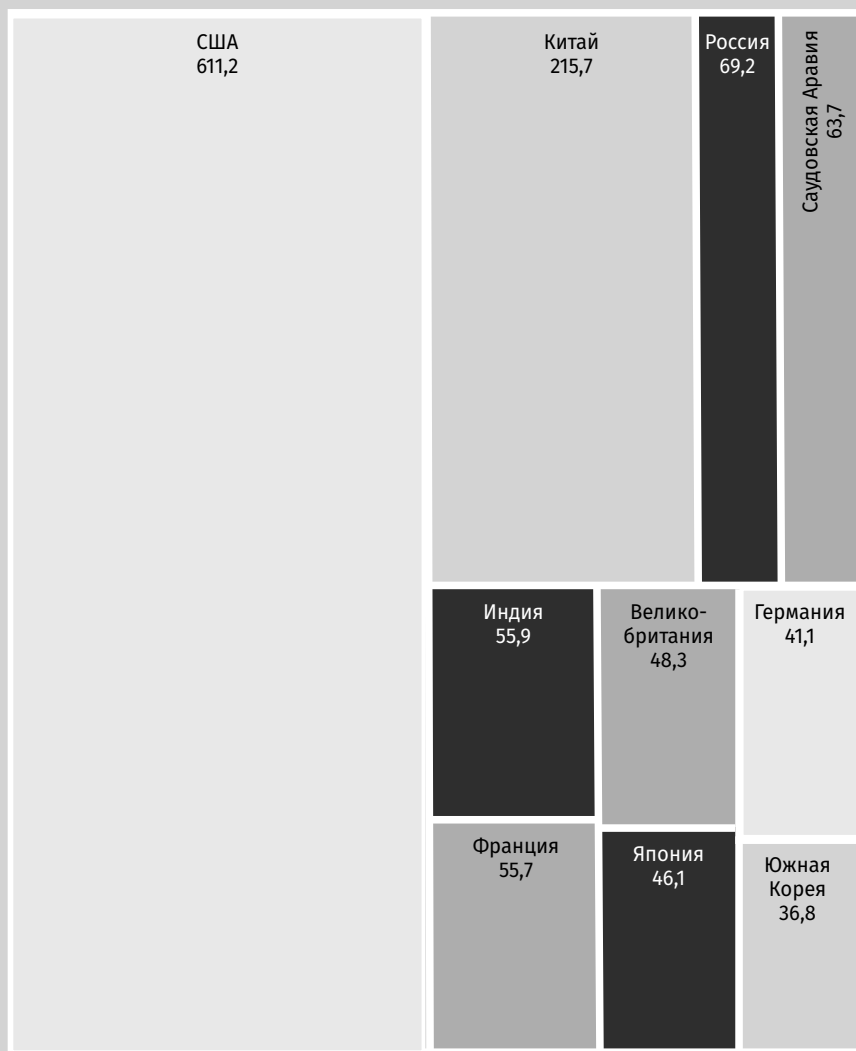


Диаграмма «водопад»

Диаграммы «водопад» позволяют отображать последовательность данных в виде положительных и отрицательных значений. Особенно полезны они при разбиении какой-либо крупной величины (например, прибыли) на компоненты.

На этой диаграмме отображены общий доход франшизы Star Wars и доля каждого фильма в нем.

Источник данных: www.statisticbrain.com

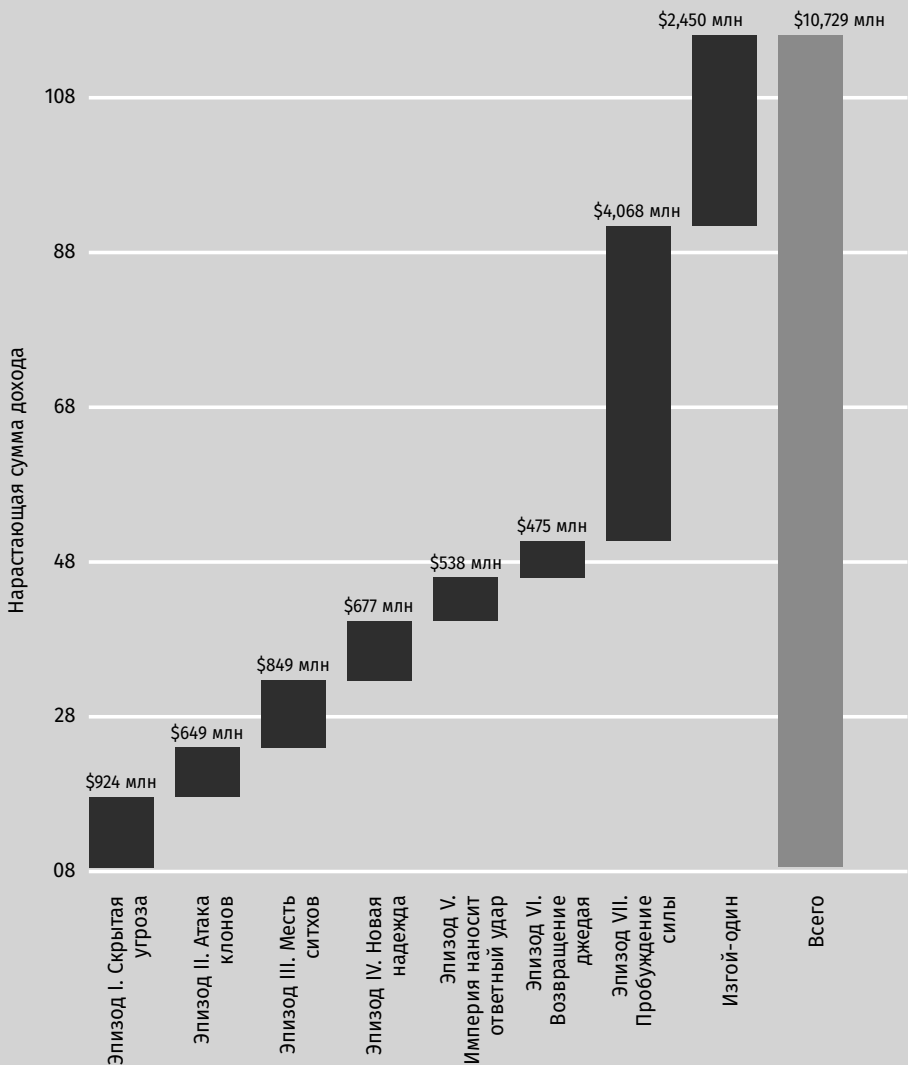
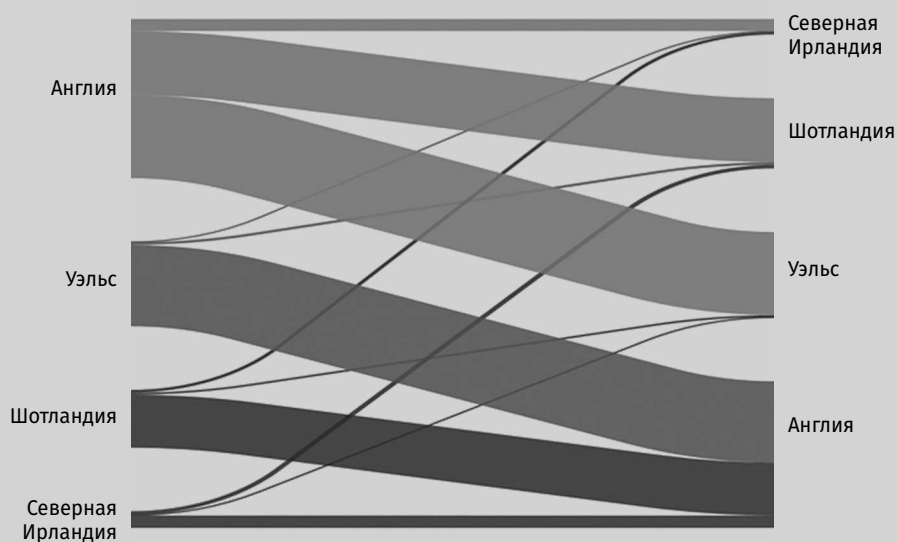


Диаграмма Сэнки

Диаграмма Сэнки отображает движение данных, используя размер и направление стрелок. Этот подход идеально подходит для визуализации любого потока данных — идет ли речь о пользователях, проходящих через воронку продаж, или о миграционных моделях.

Приведенный ниже пример демонстрирует, каковы миграционные тенденции в регионах Великобритании. Источник данных: theinformationlab.co.uk



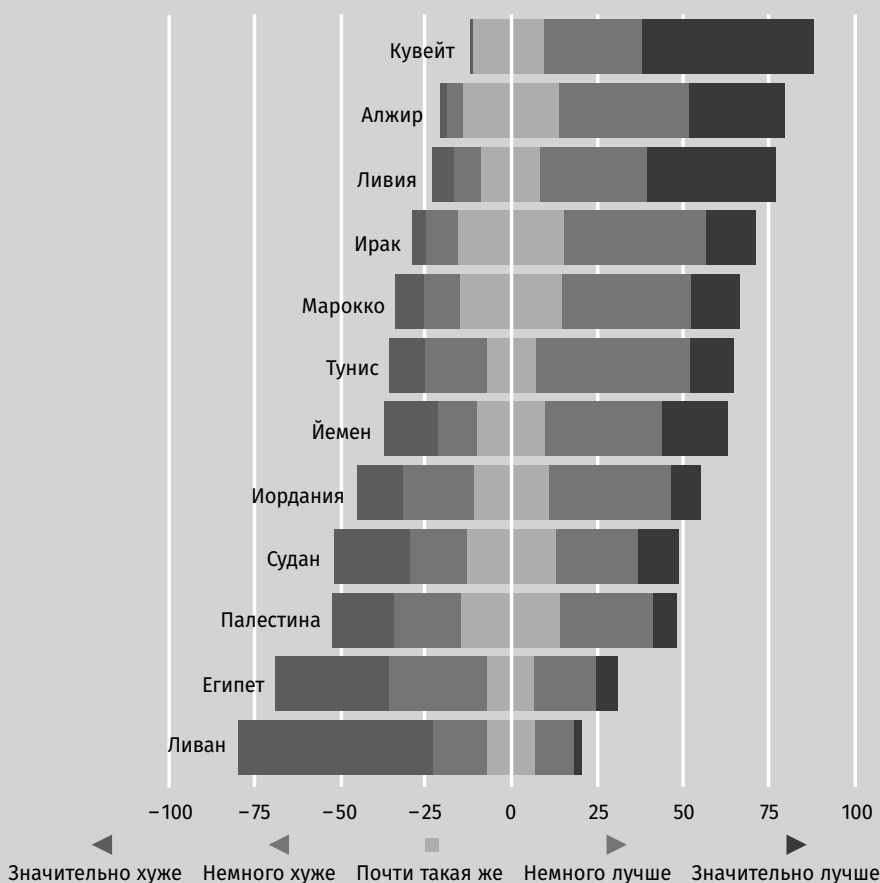
Шкала Ликерта

Это психометрические шкалы информации, представляющие собой сумму ответов на вопросы анкеты. Шкала Ликерта может использоваться только в тех случаях, когда применяется система Ликерта или аналогичная ей система (обычно для этого требуется один ответ на вопрос по скользящей шкале из пяти пунктов).

Данная шкала отражает то, какой граждане видят свою экономику в ближайшие три–пять лет.

Источник данных: www.arabbarmometer.org

В.: Как вы думаете, какова будет экономическая ситуация в вашей стране в ближайшие годы (3–5 лет) по сравнению с нынешней ситуацией?



данных, и каждый из них имеет свои преимущества и подводные камни с точки зрения изложения истории ваших данных. При использовании в неправильном контексте простая столбчатая диаграмма может оказаться такой же скучной, как и массив данных.

Заключительные размышления

Для меня стадия визуализации охватывает также предшествующий ей и последующий этапы (анализ данных и их представление) в ходе анализа и обработки данных. Визуализация может внести значительный вклад в аналитику и фундаментально способствует окончательному представлению полученных результатов заинтересованным сторонам. Так что это чрезвычайно полезный процесс, который не только помогает убедить других в силе данных, но может подарить нам самим новые идеи.

Сейчас мы переходим к финальной части — презентации наших данных. Визуализация — это последняя возможность внести изменения, перезапустить анализ и отфильтровать данные. Мы должны быть готовы представить наши результаты, прежде чем перейти к заключительному этапу работы.

Заключительный этап процесса анализа и обработки данных нацелен в первую очередь на конечных пользователей — людей, инициировавших проект. Если вы воспользуетесь советами предыдущих глав этой книги, то вам придется на каждом этапе принимать во внимание их интересы и помнить, какие результаты им нужны. Сейчас настало время перевести наши идеи в формат, который смогут понять и использовать неспециалисты по данным.

Эта глава, по сути, посвящена общению с людьми. В ней объясняется, как собранная и обработанная информация может быть уточнена и передана в виде единого согласованного сообщения заинтересованным сторонам. Если они неправильно поймут ваше сообщение, то либо вообще не совершат никаких действий, либо предпримут неправильные шаги — что будет катастрофой для любого проекта.

Именно поэтому всякий раз, когда меня спрашивают, что отличает хорошего аналитика данных от аналитика высшего класса, мой ответ всегда таков: способность передавать идеи. И эта глава расскажет, как отточить свои презентационные и коммуникативные навыки таким образом, чтобы добиться успеха.

Важность повествования

Вспомните нашу отправную точку в первой части: рассказывание историй является ключевой частью процесса науки о данных*, и *наша*

* Мой друг и наставник, Яу Тан, старший вице-президент по аналитике финансовых преступлений и управлению программами DBS Bank (Сингапур), полушутя относится к стадии представления процесса анализа данных как требующей дополнительных 80% вашего времени (помимо 100%, которые вы потратили на этапах 1–4), потому

задача — убедиться, что аудитория вовлечена в ваш рассказ. Не надейтесь на то, что, поскольку вы говорите с людьми, инициировавшими проект, они обязаны понять смысл ваших идей и результатов. Помните: они пригласили вас, потому что им был нужен эксперт, который рассказал бы о том, что говорят данные. Теперь ваш шанс сделать ответный ход! Пересмотрите каждый шаг процесса и спросите себя: «Что мне нужно делать на каждом этапе? На какие вопросы необходимо ответить? Как решить возникшие проблемы? Что показывают результаты? Что-нибудь осталось без ответа?» Задавая себе эти вопросы, вы сможете придумать историю, соответствующую вашему опыту, — в конце концов, вы рассказываете *свою* историю о том, как вы взаимодействовали с данными.

Умение передать идеи — именно то, что нужно в организациях, сотрудники которых могут быть не особенно технически подкованы. Важно как можно лучше представить данные, поскольку это может обеспечить нам преимущество перед теми, кто, жонглируя цифрами, может построить модели и написать код, но не способен объяснить, какое отношение полученные результаты имеют к компании.

Кейс: Veros Credit — представление конечному пользователю

Любой, кто брал кредит, знает, каким сложным процессом это может оказаться. Но у кредитных компаний есть и свои собственные трудности. Ведь это не благотворительные организации: они должны оценить, окупится ли кредит нового клиента или нет. А чтобы дать правильную оценку, следует учесть различные характеристики клиента (от кредитоспособности до истории занятости) и выстроить на основе этих показателей тактику, позволяющую снизить риски для компании.

Грег Попп — старший вице-президент по управлению рисками в Veros Credit, американской компании финансовых услуг, которая предоставляет кредиты людям, желающим приобрести автомобили в автосалонах. Его команда курирует проекты компании в области науки о данных и финансового инжиниринга, и одной из их центральных задач является создание модели оценки кредитоспособности для прогнозирования риска по каждому клиенту. Необходимые данные о потенциальных заемщиках определяют кредитный рейтинг, на основе которого Veros

что то, что мы делаем на данном этапе, очень важно. Я называю это «правилом Яу 80–20–80» (SuperDataScience, 2016a).

Credit будет либо утверждать, либо отклонять запросы на получение кредита. Затем эта информация уходит в автосалон, где решают, принимать или не принимать условия, предложенные Veros Credit.

Идея сводится к следующему: необходимо убедиться, что процентная ставка по одобренным кредитам и выплаты по ним покроют любые риски, связанные с кредитными и операционными расходами. Для поддержания предсказательной способности модели мониторинг оценки рисков осуществляется на ежемесячной основе. В ходе этих наблюдений параметры корректируются в соответствии с новыми данными. Такая прозорливость помогает Veros Credit выделять признаки возможных неблагоприятных результатов кредитования и позволяет отказаться от клиентов, как только их показатели начинают отклоняться от установленной нормы.

Общение с аудиторией

Это все сложная работа, но для Грега самая трудная часть — добиться согласия на управленческом и исполнительном уровнях. Ему приходится отстаивать свой проект — и, поскольку дело затрагивает многие корпоративные позиции, это следует принять во внимание при создании презентаций. Донести смысл управляемых данными моделей до инициаторов проекта часто непросто. Генеральный директор без знания предмета или даже статистики вполне может воспринять сильно различающиеся данные как свидетельство неудачи проекта. По этой причине, когда Грег готовит презентацию, он намеренно фокусируется на конечном пользователе, стараясь сделать свое сообщение понятным для тех, у кого нет опыта в области науки о данных.

«Будучи студентом колледжа, я подрабатывал репетитором по статистике, что научило меня объяснять сложные проблемы так, чтобы люди, не владевшие предметом, могли в них вникнуть. В итоге мои подопечные получали более высокие оценки — и были счастливы. Этот опыт помог мне и как аналитику данных: я использую тот же подход при объяснении своих моделей менеджерам. Если вы не можете объяснить, что делаете и почему, вам всегда будет нелегко».

(SuperDataScience, 2016b)

Я действительно считаю, что навыки презентации абсолютно необходимы аналитику данных. Не все придут в науку о данных, имея за плечами профессиональный опыт преподавания, но вы можете извлечь ценные уроки из общения с друзьями и близкими, применяя принципы, о которых узнали из этой книги. Это даст

несколько преимуществ: человек, которому вы передадите информацию, узнает о замечательном мире науки о данных, а вы будете развивать свое умение «делать сложное простым», как рекомендует миссия SuperDataScience.

Формирование позитивного отношения к данным

Наука о данных относительно новая дисциплина, следовательно, мы не можем ожидать, будто люди знают, с чем они имеют дело и насколько она полезна для бизнеса. Некоторые могут подумать, что наука о данных — причуда (помните, что ваша аудитория не читала эту книгу!) или что ее методы не могут быть применены к их собственной работе. Поэтому мы должны бороться за использование данных в будущих проектах компании и отстаивать идею, что все бизнесмены сегодня должны быть грамотными в этой области.

Если инициаторы проекта поддержали нас после нашей презентации, мы можем сделать их сторонниками науки о данных. Продемонстрируйте им ценность данных — и они расскажут другим, а те в свою очередь тоже задумаются о применении данных в работе.

Кейс: обработка сообщений

Знакомство аудитории с нашими выводами может иметь более далекоидущие последствия, чем вы думаете, — надо только приложить немного усилий. Руководители проектов не всегда потребуют от вас сделать презентацию, но я рекомендую в любом случае взять на себя инициативу — если вы сумеете доказать, что работа на основе данных может помочь компании, это принесет вам дивиденды.

Когда я работал в Sunsuper, австралийском пенсионном фонде, меня однажды попросили оптимизировать обработку поступавших от клиентов сообщений, связанных с ежегодной кампанией под названием supermatch. Сообщения обрабатывались со значительным отставанием, и потребителям приходилось долго ждать, прежде чем их запрос будет рассмотрен. Проще говоря, у компании не было возможности обработать то количество клиентских запросов, которое она получила в ходе акции.

Благодаря науке о данных я сократил срок обработки запросов с 40 дней до трех.

Это много значило для клиентов, бизнеса и доходов компании. Я также позаботился о том, чтобы моя работа не заканчивалась этими результатами. Я вернулся в операционный отдел и предложил представить в презентации свои выводы, результаты и подход, который использовал. На самом деле я *не должен* был этого делать, но, если бы я пропустил этот шаг, сотрудники отдела не поняли бы, в чем суть проделанной работы, и, возможно, не увидели бы, как наука о данных может оказаться полезной для их компании. В действительности они были так довольны итогами проекта, что продолжали обращаться ко мне с вопросами о том, насколько мы могли бы сообща повысить эффективность компании в дальнейшем.

В общем, не храните свои секреты, как жадный волшебник. Старайтесь изо всех сил объяснять своим клиентам ваш подход и то, как наука о данных может значительно улучшить их бизнес.

Как подготовить убойную презентацию

Есть множество книг о том, как сделать эффективную презентацию. С точки зрения методологии презентация результатов анализа и обработки данных не слишком отличается от презентации на любую другую тему. В следующем разделе подробно описываются методы, в эффективности которых я убедился на *своем собственном опыте* разработки презентаций. Не существует единого подхода к созданию хорошей презентации. Я разработал свой собственный стиль — он мне подходит и дает возможность чувствовать себя комфортно. Ниже приведены методы, которые, по моему мнению, могут быть полезны новобранцам в науке о данных*.

Пусть эти приемы послужат вам не только для официальных презентаций — я часто использую их, когда рассказываю новым

* Хотя эти советы могут быть полезны, не позволяйте им ограничивать ваше творчество. Не существует единого наилучшего подхода к составлению и проведению презентации. Стиль, который наиболее удобен для вас, может отличаться от того, что подходит мне.

коллегам (и даже друзьям и семье) о своей работе, или на моих онлайн-курсах. Если у вас есть повод поговорить с кем-то о науке о данных, воспользуйтесь им как возможностью попрактиковаться перед тем, как вести разговор на эту тему в более формальных условиях.

1. Подготовьте структуру

Я начинаю каждую презентацию с мозгового штурма по теме проекта. Для меня это означает запись начальной (А) и конечной (Б) точек проекта друг напротив друга на листе бумаги (часто я их просто воображаю). Пустое пространство между ними представляет мой маршрут, и я хочу показать аудитории, как я прошел путь от А до Б. В этом пространстве я записываю этапы процесса анализа и обработки данных,

Пишите, когда вы устали

Это может показаться контрпродуктивным, но мое любимое время, чтобы начать работать над презентацией, — самый конец дня, когда я устал и хочу спать. Причина этого тройная: 1) усталость ограничивает мои усилия, то есть я с меньшей вероятностью заполню страницу техническими деталями; 2) если я знаю, что у меня есть только пара часов, я скорее что-то напишу, чем если бы знал, что у меня впереди целый день; и 3) я думаю, что люди могут быть более творческими, когда они устали. Этот подход едва ли эффективен для всех, но мне он определенно помог совершить некоторые крупные прорывы в моей профессиональной деятельности. Я не знаю точно почему — возможно, некоторые части мозга, которые препятствуют творчеству, отключаются и позволяют активизироваться потокам креативности.

Конечно, я могу рекомендовать это только для первого наброска; не ожидайте проснуться на следующее утро с безупречным докладом. Однако, если вы пытаетесь преодолеть ужас перед пустой страницей, я искренне предлагаю попробовать. Даже если вы сможете использовать только 10% того, что напишете, вы все равно уже что-то сделали для презентации, за которую возьметесь на следующее утро. Мне гораздо легче редактировать написанный текст, чем писать что-то совершенно новое.

выделяя любые моменты, которые считаю особенно впечатляющими или показательными.

Такой мозговой штурм дает мне прочную структуру презентации без каких-либо дополнительных усилий с моей стороны. Внимание к форме очень важно, поскольку в ней должно быть заключено как можно больше подсказок о том, что последует дальше. Не пытайтесь оградить людей от вашей методики, думая, что ваш рассказ им будет неинтересен (или, хуже того, полагая, что *если* они не понимают, то и не будут задавать слишком много вопросов), — заставьте их принять ваш подход. Процесс обработки и анализа данных уже имеет логичную, простую в использовании структуру, так почему бы не объяснить сущность этапов и не провести слушателей по ним шаг за шагом? Когда люди чувствуют себя комфортно — ведь им ясно, о чем идет речь, — они запоминают больше информации и будут не только слышать, что вы говорите, но на самом деле *слушать* вас.

2. Соберите иллюстративные материалы

Когда во время мозгового штурма я собираюсь заполнить пустые места на бумаге перечнем этапов работы, я учитываю типы технических средств, которые смогу использовать для лучшей иллюстрации моих действий. После предыдущего этапа процесса анализа и обработки данных (визуализация) у меня уже должны быть несколько графиков, на которых я могу рисовать. В зависимости от характера проекта я иногда использую дополнительные изображения: логотипы компаний, программные интерфейсы и ассортимент продукции. Ориентируясь на знания моей аудитории, имеющие отношение к проекту, я также могу включить некоторые диаграммы или блок-схемы, которые помогут людям легче понять концепции, которые я предлагаю к обсуждению.

Помните, что есть много других средств, помимо изображений. Вы можете использовать видео- или аудиоклипы для усиления своих аргументов. Вы даже можете добавить немного юмора: шутливые мемы способны творить чудеса в презентации, особенно если речь идет о сухом предмете, — просто убедитесь, что вы можете связать юмор с вашей темой, и не переусердствуйте.

3. **Найдите свою аудиторию**

Хотя я не выступаю за преднамеренное «упрощение» или «усложнение» презентации для вашей аудитории, важно знать, кто будет находиться в зале, чтобы убедиться, что мы затронем болевые точки присутствующих. Независимо от того, с кем мы говорим — даже если с людьми, хорошо разбирающимися в науке о данных, наша презентация всегда должна быть четкой и хорошо структурированной. Аргументы и результаты следует организовать логически, и мы никогда не должны рассчитывать на априорную осведомленность слушателей.

Тем не менее, если мы хотим разбить свою аудиторию на различные группы, это может быть скорее связано с нашей собственной психологией — в частности, с тем, как мы воспринимаем своих слушателей. Самая интересная аудитория для меня — люди, управляющие компаниями: высшая исполнительная команда. Именно поэтому мы уделим им особое внимание в этом разделе.

Большинство работников точно не радуются перспективе представить свои выводы руководству компании. Представьте, что вы идете в зал заседаний и видите, как генеральный директор и четыре других руководителя пялятся на вас. Вы *должны* сказать им что-то ценное — они ждут этого. Еще хуже то, что вы знаете: они присутствовали на множестве презентаций раньше, то есть вам придется потрудиться вдвое больше, чем для любой другой аудитории.

Лучшее, что мы можем сделать, представ перед такой группой опытных слушателей, — это напомнить себе, что они всего лишь люди. Правда, они очень занятые люди, но это просто придает большую значимость четкости и последовательности изложения.

Здесь стоит упомянуть несколько других приемов, которые всегда хорошо служили мне и которые вы можете использовать для общения с руководителями.

Ноль на долларовых счетах

Руководители в основном озабочены итогами. Выясните (если вы не сделали этого раньше), сколько ваш проект сэкономит, заработает или будет стоить компании, и попытайтесь оценить его в денежном выражении. Результаты вашего проекта могут, например, увеличить

доход, или служить рекомендацией для изменения организационного или стратегического подхода, или привести к улучшению обслуживания клиентов, но не думайте, что на этом вы можете поставить точку. Изменения всегда будут *связаны с расходами* для компании, и выяснить, что это за сумма, — часть вашей работы.

Будьте конкретны. Например, если вы рассматривали пользовательский опыт, подумайте, как изменения, которые вы предложили, сохранят количество клиентов X, которые будут платить Y в течение времени Z. Попросите соответствующие отделы помочь вам с цифрами затрат и ожидаемой прибыли. Короче говоря, думайте как инвестор из реалити-шоу Dragon's Den или Shark Tank и следите за самым главным.

Будьте готовы к сложным вопросам

Нам повезло, что мы работаем в науке о данных. Хотя она быстро развивается, все еще редко можно найти ключевого аналитика данных в компании и еще реже мы сталкиваемся с кем-то на уровне правления, кто действительно в курсе передовых технических ноу-хау.

Многие практики считают, что этот недостаток знаний обезопасит их от вопросов, что часто верно для определенного типа аудитории. Но незнание руководителей только добавляет нам еще один уровень сложности — поскольку они отвечают за компанию, то будут особенно осторожны (и любознательны), когда дело касается новых подходов. Чтобы подготовиться, не переходите к обороне. Сосредоточьтесь только на фактах — это то, что нам выгодно. Данные не лгут. Они могут свидетельствовать об успешности компании или же открыть некоторые жесткие истины. В любом случае не должно быть никаких сомнений в том, что лучшее понимание данных только поможет компании. Предоставьте столько деталей, сколько им нужно, покажите преимущества и постарайтесь говорить доступным языком, не используя специальную лексику.

Вопросы и ответы: обсуждение без стресса

Вы должны подготовиться к сессии вопросов, предполагающих предоставление вами более подробной информации. Предварительное обдумывание вопросов, которые могут быть вам заданы, упростит

эту изнурительную стадию презентации. Руководители особенно склонны вдаваться в подробности, когда есть риск издержек для компании. Если вы обнаружили пробелы или аномалии в результатах анализа данных, то, естественно, захотите объяснить их в своей презентации. Но, на мой взгляд, лучше вместо этого осветить выявленные проблемы. Тогда вам будет понятно, какие вопросы задаст аудитория во время обсуждения, и вы сможете к ним адекватно подготовиться.

4. Прямо укажите проблему...

Для того чтобы донести до слушателей свои идеи, стоит обязательно напомнить им о задаче, которую перед нами поставили, — для этого можно просто вернуться к первому этапу: определить вопрос. Постарайтесь, чтобы формулировка звучала как можно более лаконично: лучше, если вы ограничитесь одним предложением.

Чтобы избежать лишнего в сообщении, представьте, что вас попросили объяснить свои действия группе шестилеток. Исключите любую информацию, касающуюся вашего метода или результатов, потому что все, что требуется от вас сейчас, — прояснить *проблему*. Только потом можно перейти к этапам, которые вы прошли, чтобы ее решить.

5. ...Затем покажите преимущества

Одно из главных правил в бизнес-маркетинге — как можно раньше показать своей аудитории преимущества того, что вы продаете. Чтобы заинтересовать слушателей презентацией, нужно объяснить, почему им следует обратить на нее внимание. После того как мы сформулировали проблему, нужно сказать, что полученные сведения помогут улучшить продажи, привлечь клиентов и т.п., — таким образом мы дадим понять заинтересованным сторонам, что аналитика данных улучшит их бизнес. А после того, как мы сообщили аудитории о возможном позитивном эффекте использования наших результатов, можно приступить к их *убедительному* представлению.

Как их увлечь

Получить поддержку аудитории не всегда просто. Многие исследования показывают, что в наш цифровой век внимание людей быстро ослабевает. Хотя мы стараемся подготовить свои выступления как можно лучше, есть некоторые факторы, которые мы не можем контролировать. Иногда мы будем общаться со слушателями в конце дня, когда они устали и хотят домой, а в другой раз можем даже встретиться с аудиторией, враждебной нашим идеям. Это означает, что во многих случаях нам придется больше работать, чтобы привлечь людей на свою сторону. Один из лучших путей достижения этой цели — взаимодействие.

В моей практике это означало либо рассмешить аудиторию, либо задать ей вопрос. В первом случае люди слышат, как окружающие смеются, и у них появляется желание вникнуть в происходящее. Если вы сможете — пошутите, но нет необходимости заставлять себя острить, если для вас это неестественно. И будьте осторожны. Безопасная шутка — шутка над собой. Я обычно говорю о том, сколь бодрящей я нахожу науку о данных, — это утверждение почти гарантированно вызывает смех. Пускай слушатели видят в вас эксцентричного ученого, и, когда вы покажете им, что к тому же можете помочь их бизнесу, они станут вашими поклонниками.

Если вы не комик от природы или если вам неудобно ломать лед таким образом, попробуйте задать вопрос, который требует поднятия руки. Например, вы можете спросить, сколько людей в комнате используют продукт или услугу, о которых вы собираетесь говорить, или сколько людей знают о концепции, которую вы планируете обсудить более подробно. Убедитесь, что ваш вопрос имеет отношение к теме презентации.

Какой способ является предпочтительным «ледоколом»? Это зависит как от вашего характера, так и от аудитории. Лично мне нравится использовать шуточный подход для более широкой аудитории. Легче заставить смеяться 100 человек, чем пять.

Статистически (я здесь использую науку о данных!) чем больше группа, тем больше людей поймут вашу шутку и рассмеются. И как только несколько человек засмеются, к ним присоединятся остальные.

Вместе с тем я считаю, что, когда вы выступаете перед руководством компании, лучше задать вопрос, потому что вы можете обращаться к людям индивидуально. Иногда мне отвечают молчанием, но я принимаю это. Если же я получаю ответ на свой вопрос, то не забываю сказать спасибо — что усиливает вовлеченность других участников обсуждения.

6. Говорите понятным языком

Исходя из собственного опыта, могу сказать, что большинство презентаций в области науки о данных сложны для понимания. Если оставить в стороне такой фактор, как неудачная подача материала, они слишком часто не структурированы, даже бессвязны и не достигают своих целей. Тому есть много объяснений: например, отсутствие у автора презентации интереса к процессу представления полученных результатов или осознанное стремление к туманному изложению предмета, чтобы лишить слушателей повода задать вопрос. Но отсутствие ясности поставит крест как на самом проекте, так и на вашей карьере в компании.

В первую очередь позаботьтесь о том, чтобы все в аудитории одинаково хорошо понимали, о чем идет речь. На презентацию придут люди, в разной степени сведущие в вашей теме, и это особенно касается такой новой дисциплины, как наука о данных. Обращайтесь к человеку, у которого, как вам кажется, знаний меньше, чем у остальных. Как и при описании перспектив проекта, постарайтесь прояснить детали максимально короткими и четкими фразами. Не только подробно расскажите, что было сделано, но и покажите, *почему* эти методы были важны для получения соответствующих результатов. Не используйте терминологию науки о данных и названия программных продуктов без краткого объяснения того, что они собой представляют или зачем нужны.

7. Расскажите свою историю

Если мы хотим продемонстрировать, насколько важна наука о данных, мы не можем ограничиваться освещением ее теоретических основ. Мы должны показать аудитории, что наука о данных имеет практическое значение. Проведите небольшую подготовительную работу, чтобы узнать об отделах, в которых работают ваши слушатели, и расскажите им о том, как методы, используемые нами, могут быть применены к их собственной работе. Это заставит их задуматься о науке о данных и ее реализации.

Краткий обзор того, как мы выстраивали свою работу, помогает раскрыть возможности науки о данных. Однако я не питаю иллюзий: разговор о некоторых областях знаний, особенно о науке о данных, с людьми, не интересующимися этой темой, может быть трудным. Этот разрыв между «нами» и «ими» становится еще больше, когда мы тараторим, просто перечисляя факты и цифры. Пошаговое описание процесса работы легче знакомит слушателей с предметом, ставя их *на наше место*, — любая аудитория скорее воспримет личную, человеческую историю, чем список данных о покупках.

Подходы мотивационных спикеров

Многие из лучших лекторов, которых я знаю, используют этот прием: делают более «человеческой» историю, которую хотят рассказать. Популярный мотивационный спикер Тони Роббинс и ученый, ставший режиссером, Рэнди Олсон, часто опираются на подходящие личные истории, чтобы продвинуть свои идеи*. Конечно, то, что Роббинс всю жизнь занимается проблемами самореализации, дает ему возможность вовсю использовать свой опыт, но можно сделать то же самое с нашими проектами по науке о данных. Все, что требуется, — это перенастроить свое мышление.

* В своей основополагающей работе «Хьюстон, у нас есть история?» (Houston, We Have a Narrative?) Олсон говорит о пользе рассказывания историй для передачи сложных научных концепций.

Вместо того чтобы размышлять о холодных, бездушных аспектах вашего проекта, подумайте о том, как вы чувствовали себя эмоционально и что с вами происходило на каждом этапе процесса. Было ли вам интересно общаться с людьми? Вас что-нибудь удивило? Вы терпели неудачи? Как ваш прошлый опыт помог вам? Что значат для вас результаты? Рассказ о том, как вы добились прорыва, зацепит вашу аудиторию.

Также помните, что, рассказывая о собственном опыте, вы не должны сосредотачиваться на положительных моментах. Люди с симпатией относятся к терпящим неудачу, потому что сами бывали в подобных ситуациях. Никто не захочет слушать историю о том, как богатый становится еще богаче. Поведайте о препятствиях, с которыми вы столкнулись, и о том, как их преодолели.

Именно такое эмоциональное выступление запомнится аудитории.

Я считаю, что именно благодаря описанному подходу моя презентация на конференции крупных суперфондов (CMSF) получилась такой, что в *Investment Magazine* появилась статья обо мне под названием «Русский специалист по данным из Sunsuper впечатлил CMSF». Если бы в ней я просто представил результаты своего анализа данных, уверен, что название публикации было бы другим. Здесь нет никакой магии. Любой может это сделать.

8. Подготовьте слайды (по желанию)

Не все презентации требуют слайдов, но лично я предпочитаю их использовать. При правильном применении слайды могут создать дополнительный уровень информации и удерживать внимание аудитории.

В ходе подготовки слайдов:

- 1) не переборщите — аудитории нужно время переварить информацию, старайтесь, чтобы слайдов было немного; показывайте не более одного слайда каждые три минуты;

- 2) ограничьте количество текста на каждом слайде; если текста будет слишком много, аудитории придется одновременно читать написанное и слушать выступающего — и в конечном итоге они будут делать и то и другое плохо!

Хотя я всегда рассматриваю презентацию в PowerPoint как средство улучшить свое выступление, хочу подчеркнуть, что слайды работают *на меня*, а не наоборот. Интровертам может не понравиться эта идея, но аудитория должна смотреть *на нас*. Мы передаем информацию — слайды лишь помогают нам.

PowerPoint и Keynote — отличные инструменты для создания презентаций. Хотя я уверен, что использую только 15% потенциала PowerPoint, все же могу составлять чрезвычайно эффективные презентации и курсы с помощью этого программного обеспечения. Большинство корпоративных организаций используют PowerPoint, а некоторые обучают своих сотрудников работать с ней. Не отказывайтесь от таких занятий, даже если вы имели дело с PowerPoint раньше, — каждый преподаватель привнесет что-то свое и покажет новые приемы создания действенного дополнительного материала.

9. Практика, практика, практика

Если ваша презентация не содержит секретной информации, я настоятельно рекомендую несколько раз прочесть ее разным людям. Тогда вы увидите, как ваш черновой проект принимают люди, (желательно) не сведущие в науке о данных. Это поможет вам отточить свою речь, сократить число терминов и обеспечить доходчивость изложения — и в итоге вы будете более уверенно себя чувствовать, когда придется представлять свои выводы «нужной» аудитории.

Вот что важно иметь в виду для «пробной» презентации:

- **Избыток слов-паразитов.** После презентации спросите слушателей, как часто вы говорили «хм» и «типа». Или лучше запишите себя на телефон и прослушайте. Постарайтесь в будущем избегать этих слов. Люди, как правило, используют слова-паразиты, когда нервничают и чувствуют, что им хочется помолчать,

хотя нужно продолжать говорить. Паузы иногда покажутся вам как докладчику невероятно длинными, но они могут оказаться полезными, ведь аудитории может понадобиться передышка в восприятии информации. Относитесь видеть в паузах не знак того, что вы отстаете, а возможность для вашей аудитории осознать услышанное.

- **Зрительный контакт.** Многие аналитики сталкиваются с этой проблемой! Хотя ситуация постепенно меняется, люди с техническим образованием в целом, как правило, в меньшей степени привыкли к публичным выступлениям. Испытывая дискомфорт в ходе своих презентаций, они могут уставиться в пол, а следовательно, потерять важнейший контакт с аудиторией. Заставьте себя смотреть на людей, к которым вы обращаетесь. Считается, что, если представить их голыми, это может помочь, правда, я никогда не заходил так далеко. Я стараюсь выбрать людей, которые кажутся наиболее открытыми и дружелюбными, и сосредоточиться на них. Дополнительный бонус для тех, кто найдет *несколько* дружелюбных лиц в разных местах зала: если вы будете смотреть то на одних слушателей, то на других, аудитория решит, что вы взаимодействуете со всеми.
- **Чрезмерная серьезность.** Наука о данных — серьезное дело, но не позволяйте этому влиять на стиль вашей презентации. Есть опасность, что вы покажетесь слишком строгим и критичным, что не расположит к вам аудиторию. Просто попытайтесь улыбнуться, даже выглядеть уязвимым. Это позволит установить связь и взаимопонимание с аудиторией, привлечь слушателей на свою сторону. Чтобы убедить людей, что вы открыты для их комментариев и мнений, расскажите им, как они могут связаться с вами и получить доступ к вашим выводам.
- **Выразительность.** Еще один побочный эффект излишней напряженности — монотонная речь. Вам не нужно быть великим актером, чтобы изменить в лучшую сторону свою манеру говорить. Поэкспериментируйте со своей интонацией и темпом, обращая внимание на моменты, когда требуется сделать

дополнительный акцент. Это очень важно, если вы не хотите, чтобы ваши слушатели заснули.

Завершение процесса

Поздравляем, мы достигли заключительной стадии процесса анализа и обработки данных! Теперь вам остается только собрать материалы, сделать глубокий вдох и представить свою блестящую работу всем участникам проекта.

Прежде чем перейти к следующей главе, задумайтесь на минуту обо всем, что вы узнали. Возможно, стоит вернуться к первой части, чтобы действительно увидеть, как далеко вы ушли. Появилось ли у вас мышление специалиста по данным? Если вы внимательно читали эту книгу, то уже должны иметь некоторое представление о своих сильных и слабых сторонах. Составьте их список на листе бумаги — этот перечень пригодится вам при чтении следующей, заключительной главы.

Ваша карьера в науке о данных

10

К 2020 г. прогнозируется, что число новых вакансий в области анализа и обработки данных увеличится на 364 000 только в США (Burning Glass Technologies и IBM, 2017). Я упоминаю эту цифру, чтобы подчеркнуть: работа в области науки о данных и впредь будет чрезвычайно востребована. Если мы вспомним Airbnb, у которой есть свои собственные университеты для обучения анализу и обработке данных (Mannes, 2017), или уволенных за ненужностью американских шахтеров, самостоятельно обучающихся программированию (Rosenblum, 2017), то нам станет очевидно, что огромное количество занимающихся технологиями (и не только) организаций сталкивается с отсутствием подходящих кандидатов на рынке. И в эпоху, когда так много рабочих мест рискуют быть ликвидированы в течение 20 лет, наука о данных должна представлять интерес для всех, кто хочет обеспечить себя гарантированной и интересной работой*.

В этой главе мы обсудим, что сделать, чтобы найти и начать подходящую работу в области науки о данных. Я расскажу о том, как увеличить свои шансы быть замеченными в этой сфере; что работодатели хотят увидеть в вашем резюме; как лучше подготовиться к собеседованию. А если у вас уже есть должность вашей мечты, мы обсудим, как вы можете упрочить свои позиции и расти.

Вхождение в профессию

Есть несколько карьерных троп в области науки о данных. В этой книге мы рассмотрели, какие функции лежат на сотруднике, выполняющем

* В отчете CrowdFlower по науке о данных за 2017 г. 88% опрошенных аналитиков данных сказали, что они либо счастливы, либо очень счастливы своим положением.

либо весь процесс анализа и обработки данных, либо отдельные его части. Требования к должности будут зависеть от учреждения и от того, есть ли там отдел по работе с данными, а также от того, какую часть бюджета организация хочет (или может) выделить команде специалистов.

Кем работать?

Теперь давайте рассмотрим, какого рода должности* могут ожидать вас. В квадратных скобках я укажу, ответственность за какой аспект процесса анализа и обработки данных лежит на работнике, занимающем ту или иную должность.

- **Бизнес-аналитик.** Такой специалист использует методы бизнес-аналитики для преобразования результатов анализа данных в графики, выводы и рекомендации. Ему всегда будут необходимы сильные презентационные навыки. [Шаги 4–5.]
- **Аналитик данных (подготовка данных).** Глава 5 показала нам, что подготовка данных является одной из самых длительных стадий процесса обработки и анализа данных, поэтому выделение этой специальности в качестве самостоятельной не должно вызывать удивления. Это позиция начального уровня, она предполагает выполнение таких задач, как очистка и структурирование данных при подготовке к анализу. [Шаг 2.]
- **Аналитик данных (моделирование)** отвечает за разработку систем и моделей, которые могут быть применены к базам данных компании. Хотя подготовка данных не всегда может входить в обязанности соответствующего сотрудника (иногда этим этапом занимается аналитик по подготовке данных), по-прежнему очень важно иметь навыки в этой области. [Шаги 2–3.]
- **Специалист по данным/расширенной аналитике/практик машинного обучения/старший научный сотрудник по данным.** Для меня это специалист в области «реальной науки

* Обратите внимание, что этот перечень не высечен в камне и должностные инструкции будут разными для каждого учреждения. Всегда необходимо внимательно прочитать описание, прежде чем подавать заявление. Представленный здесь список можно использовать в качестве руководства для написания заявления.

о данных». Профессионал, который подходит для этой работы, должен знать процесс анализа и обработки данных как свои пять пальцев, проявлять инициативу, быть ориентированным на данные, творческий подход и разбираться в программировании и анализе. Для большинства должностей также могут потребоваться навыки визуализации и презентации. [Шаги 1–5.]

- **Менеджер по анализу и обработке данных.** Это организационная должность, и поэтому не все будут считать ее значимой для развития карьеры — некоторые захотят остаться аналитиками данных. Менеджер по аналитике общается с клиентами и/или возглавляет команду, обеспечивая выделение нужных ресурсов и людей для проектов. [Шаги 1 и 5.]

Старайтесь быть гибкими. Особенно если вы только начинаете свою карьеру, не отказывайтесь от одной вакансии в пользу какой-то иной — иногда должность, к которой вы стремитесь, можно получить, поработав сначала на другой позиции в течение некоторого времени. Я научился относиться к смене ролей как к смене стадий процесса анализа и обработки данных, то есть можно начать как бизнес-аналитик или специалист по подготовке данных и трудиться, пока не достигнешь нужной позиции. Преимущество этого метода в том, что вы получите ценный опыт работы на ранних фазах процесса анализа и обработки данных, а это означает, что вы будете лучше подготовлены к выявлению и решению проблем, с которыми можно столкнуться на последующих этапах.

Скрытые должности

Обратите внимание, что вакансии, связанные с наукой о данных, не ограничиваются перечисленными выше. Есть ряд позиций, которые, как выяснилось, отлично подходят для людей, только что окончивших университет (или другие учебные заведения). Они хотят попробовать что-то новое, прежде чем с головой погрузиться в свою профессию.

Эти должности немного похожи на должности помощников юристов в юридической фирме — будьте готовы работать на кого-то и, возможно, выполнять задачи, которые вы могли бы счесть несущественными. Но если вы трудитесь в большой компании, такого рода деятельность может открыть перед вами широкие перспективы и стать отличным шагом к более высокой позиции (и, как вы узнаете позже в этой главе, знакомство с коллегами и работа в команде — ключ к успеху любого специалиста по данным).

Еще один путь, который могут выбрать читатели, более уверенные в своих умениях и знаниях, — консалтинг. Пусть это слово не смущает вас: если в других отраслях люди становятся консультантами только после того, как они приобрели многолетний богатый опыт в своей области, консультантом по науке о данных может стать и новичок (как и я) — в таком случае он будет фактически являться советником руководителя компании, желающего знать, как читать свои данные. Консультанты могут даже привлекаться для принятия важных решений и разработки политики в отношении науки о данных в организациях. В качестве консультанта вы не только получите возможность участвовать во всем процессе анализа и обработки данных, но и, если правильно разыграете свои карты, больше узнаете о различных отраслях. Это даст вам весомое преимущество перед конкурентами, так как через пару лет работы вы, скорее всего, будете представлять, в какой сфере хотели бы развивать свою карьеру.

Консультирование — та деятельность, которая позволит вам избежать преждевременной узкой специализации. Это особенно верно для небольших консалтинговых аналитических фирм, где количество сотрудников невелико и поэтому каждый должен стать своего рода многофункциональным «швейцарским армейским ножом» для обработки данных (SuperDataScience, 2017b). Такие фирмы предоставляют отличные возможности для начинающих.

Да, консалтинг означает долгие часы и тяжелый труд, и по этой причине такая работа может быть неидеальна для тех, у кого есть семья. Однако если вы способны и готовы пожертвовать своей личной жизнью на несколько месяцев ради начала вашего проекта, то в итоге получите важный опыт, узнаете особенности отрасли и то, как в нее вписывается наука о данных. По сути, работа в качестве консультанта поможет вам определить, чем вы хотите заниматься в будущем.

В какой области?

Поскольку специалисты по данным пользуются таким высоким спросом, читатели наверняка захотят узнать, как разработать стратегию, чтобы сосредоточиться на конкретной области. Многие компании активно ищут аналитиков данных, и для людей, которые еще не выбрали для себя сферу деятельности, было бы логично пойти туда, где ожидается рост спроса. В докладе Burning Glass Technologies и IBM за 2017 г. рассматривается спрос на специалистов по работе с данными в шести ключевых секторах (профессиональные услуги; финансы и страхование; производство; информация; здравоохранение и социальная помощь; розничная торговля). В документе отмечается, что отрасль профессиональных услуг (которая, как правило, включает в себя в том числе консультации по вопросам управления, юриспруденции и медицины) дает львиную долю вакансий, а финансовая и страховая сферы занимают второе место.

Вывод, который мы можем сделать: если вы хотите получить наибольший шанс по крайней мере на этапе интервью, обдумайте возможность применить свои знания в области профессиональных услуг или в финансовой и страховой отраслях, где предлагается больше вакансий.

Еще один хороший способ найти работу аналитика данных — просто следить за новостями. Как говорилось в главе 2, данные и искусственный интеллект являются горячими темами на текущий момент, и внимание к ним со стороны СМИ только продолжит расти, поэтому

будьте в курсе и заведите папку с вырезками статей о тех областях, которые вас интересуют.

Даже успешные ученые должны продолжать следить за развитием событий. Дэн Шиблер, аналитик данных в True Motion, говорит:

«Это век информации. Когда мне становится любопытно, как наука о данных развивается в другой области, мне достаточно проглядеть некоторые исследовательские работы, чтобы узнать, как идут дела. Я заинтересовался новейшими разработками в области распознающих изображения сверточных нейросетей и их связью с нейронаукой. Я нашел в Университете Брауна профессора, который провел много исследований по этому вопросу, и в свободное время тружусь вместе с ним. Здорово, что я могу работать в передовой сфере и заниматься вещами, полностью противоположными тому, что я делаю в TrueMotion».

(SuperDataScience, 2017a)

Помните: для того чтобы преуспеть и чтобы ваши идеи не устарели, вы должны читать специальную литературу и постоянно развиваться. Делайте как Шиблер и отправляйтесь на охоту за проектами, которые вы сочтете увлекательными.

Что еще?

Вы уже сделали первый шаг в деле изучения науки о данных, прочитав эту книгу. Но здесь важно подчеркнуть, что это может быть только *первым* шагом. Одной книги никогда не бывает достаточно — да и не должно быть!

Вот что вы должны сделать, прежде чем подать свое первое заявление о приеме на работу.

1. Дайте себе (больше) времени, чтобы узнать свое ремесло

Осознав, насколько высок спрос на науку о данных, университеты начали обучение студентов по таким специальностям, как «магистр наук в области прогностической аналитики (онлайн-курс)» (Северо-Западный университет), «бизнес-аналитика и социальные сети»

(Университет Брунеля), «вычислительная статистика и машинное обучение» (Университетский колледж Лондона).

Если честно, я считаю, что нет большого смысла в получении таких дипломов. Помимо финансовых издержек, связанных с обучением, есть еще одна проблема: эта дисциплина *очень быстро* стареет. Разработка новой учебной программы в университете может занять более шести месяцев. К тому времени, когда курс пройдет через все бюрократические препоны, сама наука, вероятно, уже уйдет дальше*.

Более того, практически всю информацию, необходимую для овладения наукой о данных, можно найти в интернете. Многие практикующие аналитики данных (включая меня) ведут онлайн-курсы, и с их помощью вы можете научиться не только теории, но и практическому применению своих знаний. Это, на мой взгляд, сейчас более разумный вариант. Но важно не терять бдительность при выборе преподавателя: в наш демократический век интернета любой и каждый может создать собственный курс. А вот найдя правильного руководителя, вы получите открытый доступ к информации, которая будут обновляться по мере того, как новые методы станут доступными, а старые практики уйдут в небытие.

2. Важно не то, что вы знаете, а то, как применяете знания

Как вы думаете, что было самым мощным инструментом в начале моей карьеры? Программирование? Tableau? Реляционные базы данных? Выборка Томпсона?

Вы удивитесь, узнав, что я имею в виду программу PowerPoint, но именно так и было. Этот простой инструмент изменил для меня всё. Из всех сложных программных пакетов и алгоритмов, которые я использовал, PowerPoint оказала на меня наибольшее влияние, и в значительной степени благодаря ей я получил свою самую заметную должность. Как я уже сказал в главе 9, есть много талантливых жонглеров числами, но аналитики данных, способные эффективно донести свои идеи, встречаются редко.

* Обратите внимание, как увеличение числа университетских степеней и курсов в области науки о данных показывает, что мир начинает признавать науку о данных в качестве самостоятельной дисциплины.

Речь здесь не о том, что PowerPoint — палочка-выручалочка для всех специалистов по данным, но не стоит недооценивать что-либо из-за простоты. Просто-напросто я понял, как сделать так, чтобы мои знания и сильные стороны работали *на меня*. Еще учась в университете, я знал, что презентации в PowerPoint особенно хорошо даются мне.

Это пример того, как я применил свои знания для своего профессионального роста. Ваши знания, скорее всего, будут лежать в иной плоскости — найдите время, чтобы понять в какой, и запишите несколько идей. В «Будущем разума» Митио Каку цитирует нейробиолога Ричарда Дэвидсона, чтобы показать, насколько важно расширять свое представление о навыках и способности к успеху:

«Ваши оценки в школе, ваши баллы, полученные на выпускных экзаменах, менее значимы для жизненного успеха, чем ваши способности сотрудничать, управлять эмоциями, отложить удовольствие и фокусировать внимание. Эти навыки гораздо важнее...»

(Ричард Дэвидсон, цит. по книге Каку)

Вы не сумма ваших дипломов. Склонность к аналитике сослужит вам хорошую службу, но общительность и интерес к этике — дополнительные важные факторы, которые часто игнорируют даже наиболее проницательные аналитики данных. Когда у вас будет время, подумайте, как эти навыки или знания могут быть применены к любому из шагов в процессе анализа данных.

3. Делитесь — и делитесь по-братски

Наука о данных — не обычная дисциплина: те, кто ею занимается, стремятся поделиться результатами своих исследований и разработок друг с другом. Давайте извлечем из этого преимущества. В интернете есть много контента с открытым исходным кодом — им можно воспользоваться, когда будет нужно попрактиковать навыки анализа и обработки данных. Как я уже отметил в главе 3, работа над тренировочными заданиями и применение новых алгоритмов к реальным массивам данных обязательны, если вы стремитесь поспевать за развитием науки о данных и хотите, чтобы ваши идеи всегда оставались

актуальными. Так что присоединяйтесь к общему разговору и поделитесь вашими данными и открытиями с коллегами*.

Если вы еще не готовы внести свой вклад, есть много краудсорсинговых проектов, где требуются специалисты для работы с большими массивами данных, способные оптимизировать их алгоритмы. Поиск в Google по фразе «проекты краудсорсинга данных» (data crowdsourcing projects) поможет вам найти их, а дополнительным преимуществом станет то, что вы сможете работать над многими из этих проектов не выходя из собственной комнаты.

Кэролайн Макколл, ассоциированный партнер FutureYou (небольшой австралийской фирмы, занимающейся консалтингом в области науки о данных), особо подчеркивает, что обмен знаниями внутри сообщества — эффективный и благодаря интернету практически бесплатный способ стать услышанным и повысить свою профессиональную репутацию. Самые успешные аналитики данных, которых я знаю, ведут блоги, записывают влоги (видеоблоги), делятся кодами на Github (сервис исходных кодов хостинга) и подобных платформах, отвечают на вопросы пользователей на таких сайтах, как StackExchange и Quora, работают фрилансерами в Upwork (крупнейшая сеть фрилансеров) и выступают на конференциях. Если вы уверены в себе или хотите сосредоточиться на определенной нише, блоги — эффективный путь к совмещению исследовательской деятельности и карьерного роста. Не беспокойтесь о том, что поначалу можете ошибиться, — вы быстро научитесь у членов сообщества или же попросите своего наставника (см. ниже) предварительно ознакомиться с вашим материалом. Будьте скромны и открыты для критики.

Кейс: поиск наставников

У меня есть несколько наставников, и не побоюсь сказать, что именно им я обязан своим успехом на поприще науки о данных. Наставники дадут вам необходимые знания и советы и смогут удержать вас от совершения ошибок, которые,

* Если вы ищете место для общения с аналитиками данных, то я хотел бы пригласить вас взглянуть на SuperDataScience, социальную платформу, полностью посвященную науке о данных.

возможно, они сделали сами, когда начинали. Я активно призываю всех вас, на каком бы этапе карьеры вы ни находились, найти наставника. Это не всегда просто: вы не можете обратиться по электронной почте к кому-то, кого вы не знаете, с просьбой стать вашим наставником — нужно сначала наладить отношения.

Виталий Долгов — один из лучших в мире независимых консультантов по развитию бизнеса, на мой взгляд. Мы познакомились в 2012 г., когда его сосед по дому попросил меня и Артема Владимирова (с которым мы встречались в главе 6) помочь передвинуть мебель. Та встреча была совершенно не связана с работой, но даже тогда я видел, что Виталий — человек, который может вдохновить меня и помочь расти. Я позаботился о том, чтобы пересечься с ним позже и хотя бы обсудить кое-какие книги. Так мы познакомились. Жизнь иногда дарит нам неожиданные возможности, и нужно сразу хватать быка за рога. Вы никогда не знаете, когда и где вы можете встретить своего наставника.

Говорят, что удача приходит, когда подготовка встречается с возможностью. Так как же подготовиться к встрече с будущим наставником? Я спросил Ричарда Хопкинса, директора PricewaterhouseCoopers и одного из моих наставников (у него есть несколько своих!), что они об этом думают. В ответ я услышал, что потенциальные подопечные должны быть в первую очередь склонны к самоанализу и, вместо того чтобы ожидать помощи от наставника, думать о том, чем они сами могут поделиться. Это может показаться невыполнимой задачей: если вы только делаете первые шаги в области науки о данных, то как можете помочь экспертам? Хотя вы пока не способны оказать существенное содействие, у вас все равно могут быть полезные навыки или опыт. Когда найдете потенциального наставника, подумайте о том, как вы можете пополнить его знания, и обратитесь к нему с соответствующим предложением — тогда ваш возможный ментор будет думать о ваших отношениях как об улице с двусторонним движением, а не как о трате собственных ресурсов и времени.

Еще один безопасный способ добиться известности в отрасли — самому стать куратором контента. Это означает нечто не более сложное, чем создание профиля в социальных сетях и ссылок на интересующую вас информацию. Правда, чтобы быть услышанным, куратору контента приходится пройти более сложный путь, чем создателю контента.

Конкуренция

В такой стратегически важной сфере, как наука о данных, не может не быть механизмов государственной и частной поддержки новых проектов. На мой взгляд, веб-ресурс Kaggle, специализирующийся на соревнованиях по предсказательному моделированию и аналитике, является лучшим помощником для продвинутых специалистов по данным, желающих играть более заметную роль в профессиональном сообществе. Возможность получить миллионы из призового фонда следует рассматривать только как бонус.

Даже если вы еще не готовы конкурировать, можете найти на платформе много интересных реальных массивов данных. Я рекомендую по крайней мере посетить Kaggle или другие подобные сайты — быстрый просмотр их веб-страниц поможет рождению у вас собственных идей.

4. Создайте сеть

Кэролайн Макколл помогает людям, обучавшимся науке о данных, найти работу в различных компаниях — от стартапов до известных корпораций. Она говорит, что в сфере, для многих являющейся чем-то абсолютно новым, создание прочной сети контактов очень важно.

FutureYou в течение года проводит семинары по созданию историй, связанных с данными. В ходе мастер-классов под названием «Расскажи мне историю» представители целого ряда отраслей обсуждают, как повлиять на принятие решений в компаниях и как добиться взаимопонимания и доверия. Эти встречи и мероприятия FutureYou — первое явление такого рода на рынке аналитики данных; когда я начинал, подобную помощь никто не оказывал.

Это всего лишь *одна* возможность наладить контакты — есть много других (хакатоны* или встречи и конференции). Идите и принимайте участие — нет ничего плохого в том, чтобы вначале посидеть в по-

* Мероприятия, на которых программисты и аналитики данных сотрудничают в технологических проектах.

следнем ряду и просто послушать коллег. Опять же, вы сможете занять место в сообществе. Даже если вы не выступаете с докладом, все равно найдутся возможности установить контакты (например, во время обеда, ужина и кофе-брейков). И даже если вам не хватит смелости заговорить с интересующим вас человеком, у вас появится повод связаться с ним в интернете и сказать, что вы были заинтригованы его идеями, — это гораздо эффективнее, чем просто «холодная» попытка сблизиться.

Завязать связи: десять советов Макколл

Кэролайн Макколл нашла подходящую работу многим подававшим надежды специалистам — новичкам в области аналитики данных. Она говорит о роли сетевого взаимодействия и выделяет в нем десять важнейших шагов (SuperDataScience, 2017b):

- 1.** Выберите специализацию или отрасль, которая вам больше всего нравится (тогда вы скорее найдете именно ту работу, которая делает вас счастливым).
- 2.** Найдите трех наставников, которые помогут вам сориентироваться и начать карьеру в науке о данных (наставники могут давать советы о том, где вы должны искать работу, помочь вам определить свои сильные и слабые стороны и связаться с нужными людьми).
- 3.** Сотрудничайте с теми, кто может усовершенствовать ваши навыки (наиболее успешные аналитики данных готовы проявить инициативу — запустить свои собственные проекты вместе со специалистами, способными развить их навыки).
- 4.** Организуйте и проведите мероприятие, даже если оно рассчитано только на пятерых (тем самым вы докажете, что достаточно дисциплинированы и хотите учиться у других, а также готовы что-то отдавать).
- 5.** Участвуйте в подготовке подкастов, посвященных науке о данных (это простой способ сделать так, чтобы ваш голос услышали и ваши идеи распространились в международном сообществе коллег).

- 6.** Заполните профиль на LinkedIn и назовите в нем проекты, над которыми вы работали (LinkedIn помогает продемонстрировать свои возможности и предложить услуги компаниям, нуждающимся в аналитиках данных, так что чем больше потенциальные работодатели узнают о ваших проектах, тем лучше).
- 7.** Ставьте перед собой еженедельные достижимые цели (если вы сомневаетесь, начните с малого и для простоты перечислите по порядку свои задачи по налаживанию профессиональных контактов).
- 8.** Посещайте встречи и конференции в интересующей вас области (будьте в курсе и оставайтесь активными — науку о данных представляет процветающее сообщество).
- 9.** Ведите блог (статьи, если они хорошо написаны и умело распространялись, могут обеспечить вам статус эксперта в своей области).
- 10.** Стремитесь познакомиться с другими людьми, даже если это не принесет вам прямой выгоды (вы воспользовались советами и помощью других и теперь должны внести свой вклад).

Заявление о приеме на работу

К настоящему времени вам должно быть ясно, что, поскольку наука о данных еще не достигла своей зрелости как дисциплина, многие компании только начинают разрабатывать стабильную политику использования анализа данных. Если со студенческой скамьи вы еще не работали в профессии, то, вероятно, уже не раз попадали в замкнутый круг: большинство вакансий требуют многолетнего опыта, но вы не можете получить свою первую должность, если у вас нет этого опыта.

Здесь я бы сказал, что, даже если в объявлении о вакансии говорится о пяти-десяти годах опыта, вы должны все равно попытаться стать претендентом. Описать все функции аналитиков данных на сегодняшний день довольно сложно, и соответствующие вакансии будут оставаться открытыми в среднем на пять дней дольше других. Я бы

учел это и предположил, что под опытом имеется в виду скорее период, в течение которого вы изучали науку о данных (либо самостоятельно, либо в учреждении) и *практиковали* ее*.

Получить реальный опыт очень просто — достаточно вернуться к массивам данных, которые я перечислил во введении к второй части, и пройти весь процесс анализа данных, описанный в этой книге. Платформы Kaggle и SuperDataScience также регулярно публикуют задачи, которыми вы можете заняться. Не бойтесь упомянуть этот опыт в своем резюме и использовать его как возможность рассказать о своих результатах. Наниматели оценят, что вы активно работаете с данными.

По сути, компании хотят найти кого-то, кто: 1) разбирается в данных; 2) может донести свои идеи и 3) поможет им сохранить конкурентоспособность. В своем заявлении убедите работодателей, что вы тот, кого они ищут, и что применение науки о данных благотворно скажется на их основных показателях.

Творить добро

Если вы не прошли первый этап отбора, подумайте о том, чтобы заняться волонтерством**. Существует ряд организаций, таких как DataKind и DrivenData, которые руководят проектами и проводят соревнования по анализу данных. Поучаствовать в их деятельности — совсем неплохо для начала.

* Даже Артем Владимиров, с которым мы встречались в главе 6, вышел на профессиональное поле с относительно небольшим опытом в этой дисциплине. Он присоединился к Deloitte после получения диплома бухгалтера, даже не умея программировать. В итоге Артем сделал впечатляющую карьеру в области науки о данных: он решает основные аналитические задачи и выступает с презентациями по всему миру.

** Национальный фонд Великобритании по науке, технике и искусству (NESTA, www.nesta.org.uk) поможет тем, кому нужна дополнительная информация о благотворительной деятельности в области науки о данных. См. особенно Baeck (2015) и Symons (2016), чтобы начать размышлять о том, какую пользу может принести использование данных.

Ряду благотворительных организаций требуются добровольцы, готовые поработать с данными, чтобы генерировать информацию, необходимую для получения финансирования этими организациями, для усиления согласованности их действий или просто для распространения сведений об их работе. Поспрашивайте, поищите в интернете, и, если есть благотворительная организация, которой вы хотите помочь, тогда, возможно, стоит прийти к ним с собственной идеей.

Подготовка к собеседованию

Вы прошли первый этап отбора, и теперь компания приглашает вас на собеседование. Как подготовиться? Мой лучший совет: вместо того чтобы тратить силы на запоминание любимых алгоритмов и программ, узнайте все, что сможете, об отрасли, в которой вы предположительно будете работать. В конце концов, вполне вероятно, что на собеседовании вас будет расспрашивать не только аналитик данных, но и руководитель или менеджер, не очень хорошо представляющий специфику вашей деятельности, поэтому будьте готовы как к техническим, так и к нетехническим вопросам. Широта кругозора имеет решающее значение. Точно так же, как вы должны подумать о том, чем можете поделиться с вашим наставником, вы должны прикинуть, как сможете *помочь* этой компании. Каким образом наука о данных способна встряхнуть отрасль? Какие болевые точки она может устранить? Что делают конкуренты компании и как вы можете улучшить ситуацию с помощью науки о данных? Решения каких задач может ожидать от вас компания?

Когда вы доберетесь до самого собеседования, проявите энтузиазм в отношении как данных, так и того, что делает компания. Покажите, что вы можете сделать (если можно привести примеры работ, упомянутых в вашем резюме, — тем лучше), и будьте готовы к вопросам.

Вопросы Ферми

Вы можете ожидать, что услышите заурядные вопросы, которые задают всем («Чем вы хотите заниматься через десять лет? Можете ли вы привести пример того, как работаете в команде? Что вы считаете своим самым большим успехом?»), а иногда вы должны быть во всеоружии, чтобы обсуждать профессиональные подходы, знание которых обязательно для кандидата на вакантную должность. Хотя я не могу помочь вам подготовиться к этим вопросам, но могу освежить вас о вопросах Ферми*. Они предназначены для проверки вашей логики, и их часто используют в собеседованиях с кандидатами на должности, предполагающие наличие аналитических способностей.

Первое, что вы должны понять: от вас не ждут *правильного* ответа — вопросы нужны, чтобы узнать, можете ли вы логически прийти к разумному выводу.

Предположим, что у вас спросили: «Сколько красных автомобилей в настоящее время ездят в Австралии?» Подумайте, как бы вы сами ответили на этот вопрос, а затем прочитайте мой ответ ниже. Помните: наниматели хотят видеть доказательства наличия логического мышления. Вот мой ответ:

«Это очень интересный вопрос. Без каких-либо цифр я не могу дать точный ответ, но могу предложить оценку, основанную на предположениях. Во-первых, сосредоточимся только на частных автомобилях, поскольку машины для коммерческих перевозок редко бывают красными. В Австралии около 24 млн жителей, допустим, что в среднем семья состоит из четырех человек. Это означает, что в стране проживает около 6 млн семей. Основываясь на том, что я знаю из личного опыта, семьи имеют от одного до трех автомобилей, поэтому давайте усредним это до двух машин на семью. Это в среднем составит 12 млн автомобилей в Австралии.

* Названы в честь Энрико Ферми. Во время первого испытания ядерного оружия в 1945 г. он оценил мощность атомной бомбы, основываясь на том, как далеко разлетелись клочки бумаги, подброшенные им вверх во время взрыва. Названное Ферми значение мощности оказалось близко к действительному.

Теперь цвета. Есть много разных цветов, но в качестве отправной точки возьмем семь цветов радуги. Добавим белый и черный, что даст нам девять цветов. Но мы также можем с уверенностью сказать, что в Австралии белые автомобили встречаются чаще всего, поскольку здесь жарко, а белая поверхность поглощает меньше всего солнечного света, поэтому я бы посчитал белый цвет дважды, учитывая его популярность. Это дает нам десять цветов. Поскольку красный является одним из них, мы можем разделить все количество автомобилей в Австралии на десять для того, чтобы узнать, сколько в стране красных автомобилей. Итак, 12 млн разделить на десять — будет 1,2 млн. То есть в настоящее время в Австралии может быть 1,2 млн красных автомобилей».

Как вы видите из моего ответа, главным были не цифры, а то, что я проявил способность ловить на лету и логически решать проблему. Мне задали этот вопрос на собеседовании, когда я претендовал на свою нынешнюю должность в Sunsuper. Что особенно любопытно, у меня отсутствовал многолетний опыт, на который они рассчитывали (они хотели шесть лет, а у меня было только три года), и это место предполагало зарплату вдвое больше моей тогдашней зарплаты. Как вы можете себе представить, я почти не имел шансов быть принятым. Но я был уверен в себе. Этот вопрос прозвучал в самом конце собеседования, и когда я ответил тирадой, то по выражению лица менеджера по найму увидел: меня взяли на работу.

Расспросите работодателей

Недаром говорят: вы задаете им столько же вопросов, сколько они задают вам. Предварительно соберите сведения о компании и на собеседовании постарайтесь разузнать у работодателей, чем конкретно она занимается. Это существенно по двум причинам.

Прежде всего компании ценят соискателей, которые нашли время поинтересоваться тем, что им предстоит делать. Это признак серьезного

отношения к работе. И, конечно, очень плохо, если претендент не в состоянии ответить на вопрос о том, какая информация о компании ему нужна. Всегда имейте наготове по крайней мере два вопроса.

Вторая причина гораздо важнее. Собеседование — возможность узнать о вашей будущей должности, культуре компании, команде, с которой вы будете работать, возможностях для роста и т. д. Слишком много людей в настоящее время застряли на позициях, которые они ненавидят, работая в компаниях, которые они не уважают. Я здесь, чтобы сказать вам: вы заслуживаете лучшего! У вас одна жизнь, и, присоединившись к компании, вы отдаете ей свой самый драгоценный ресурс — время. Я умоляю вас, выбирайте мудро, выбирайте сердцем, кому вы отдадите свое время. Собеседование нужно не только для компании, чтобы узнать, соответствуете ли вы ее требованиям, но и для вас — чтобы понять, подходит ли эта компания вам.

Развитие карьеры в компании

Вы получили должность. Прежде всего — поздравляю! Теперь можно начинать расправлять крылья в компании. Чтобы помочь в этом, я хочу рассказать, как создал себе имя в науке о данных.

Я считаю, что мне очень повезло получить работу в Deloitte сразу после окончания университета. Но даже несмотря на то, что я был выпускником, я не позволил отсутствию опыта встать на пути моего карьерного роста — и убедил окружающих, что ко мне можно обращаться по всем вопросам, связанным с данными. Люди *хотели*, чтобы я участвовал в их проектах. И часть этого успеха была связана с тем, что я активно объяснял, как данные могут быть использованы в работе компании.

Помогая людям увидеть всеобъемлющую пользу данных, вы обретете поддержку сторонников использования данных (и, следовательно, вашей деятельности) в организации (см. главу 9). Мыслите как аналитик данных и *поделитесь результатами своих исследований*. Организуйте встречи с коллегами, чтобы продемонстрировать итоги своей работы, — и помните в первую очередь о тех, кто помогал вам в этом конкретном проекте.

В Deloitte один из проектов, над которым я трудился, включал разработку информационной панели системы управления (MIS) для розничной аптеки. Сам процесс был прост, но усовершенствования, которые внесла наша команда, в конечном итоге сэкономили компании \$19 млн. О подобных проектах вы, возможно, захотите поведать коллегам — чтобы рассказать о полученных результатах, а также чтобы люди знали: вы открыты для обсуждения того, что еще можно улучшить с помощью данных. Начало разговора — ключ к созданию аналитической культуры.

Еще один способ, которым вы можете развивать свою карьеру в компании, — проявить активность в поиске новых путей оптимизации ее деятельности. Данные генерируют идеи. Довольно часто компании не знают, что делать с аналитиками данных после того, как те выполнили поставленную задачу. Это может быть чревато тем, что ваш контракт не будет продлен после завершения проекта. Или же вы начнете скучать — что плохо для всех. Боритесь за свою позицию (вы ее заслужили!); сформулируйте проблему, обоснуйте необходимость ее решения с помощью данных и спросите своего линейного менеджера, можете ли вы заняться ею в менее напряженные периоды. Ваше отношение к работе, скорее всего, будет оценено по достоинству.

Наконец, я вновь и вновь говорю о том, насколько важно как можно раньше в своей карьере овладеть искусством презентации.

Что делать, если компания не реагирует

Не всегда все идет гладко. Если компания, несмотря на ваши усилия, не проявляет интереса к данным или если ваши просьбы начать новый проект остаются без ответа, вы можете рассмотреть возможность применения своих сил в другом месте. Я не говорю, что вы должны в этом случае бросить работу, но важно быть внимательным к тому, как вы проводите свое время и в какой области желаете расти. Я предполагаю, что вы читаете эту книгу, потому что хотите построить захватывающую карьеру в области науки о данных, и если вы продолжите работать в компании, которая не видит преимуществ использования данных, то в конечном итоге будете разочарованы.

Если бы я не разработал свой уникальный стиль презентации, то был бы просто еще одним хорошим аналитиком данных. И хотя такие специалисты, несомненно, важны, их много вокруг. Примите это как стимул к тщательному усвоению основных моментов главы 9. Когда вы представляете свои результаты, вы становитесь связующим звеном между наукой о данных и лицами, принимающими решения на основе ваших выводов. Это неизбежно заставит вас более активно и убедительно отстаивать то, что вы делаете.

Вы достигли конца книги. Молодцы, отлично потрудились; было приятно совершить с вами это путешествие. Прежде чем отправиться дальше, я бы рекомендовал не спешить и переварить все, что вы узнали. Подумайте о том, что вас действительно интересует, в каких сферах проявляются ваши сильные стороны, а в каких — слабые. Тогда вы сможете решить, будет ли вам полезно пройти дополнительный курс, — есть много ресурсов, доступных на superdatascience.com, которые помогут вам продолжить совершенствование своих навыков. Желаю успешных следующих шагов, и не забудьте время от времени давать мне знать, как у вас дела.

До новых встреч, удачи в анализе данных!

Благодарности

Я хотел бы поблагодарить моего отца Александра Еременко, чья любовь и забота сформировали меня таким, каким являюсь сегодня. Под строгим руководством отца я понял, как могу ухватиться за шансы, которые дает жизнь. Спасибо моей щедрой, отзывчивой маме Елене Еременко за то, что она всегда давала мне возможность высказать свои безумные идеи и вдохновляла меня и моих братьев на постижение более широкого мира — через музыку, язык, танцы и многое другое. Если бы не ее мудрый совет, я бы никогда не поехал в Австралию.

Спасибо моему брату Марку Еременко за то, что он всегда верил в меня, и за его непоколебимость. Его бесстрашие в принятии вызовов, которые ему бросает жизнь, подпитывает множество важных решений, которые я принимаю. Спасибо моему мудрому не по годам брату Илье Еременко за его впечатляющие бизнес-идеи и продуманные начинания. Я уверен, что слава и удача скоро постучатся и в его дверь.

Спасибо моей бабушке Валентине, тете Наташе и двоюродному брату Юре за их бесконечную любовь и заботу. И спасибо семействам Танакович и Сворен, включая моих братьев Адама и Дэвида, за все самые лучшие моменты, которые у нас были.

Спасибо моим студентам и тысячам людей, которые слушают подкаст SuperDataScience. Это вдохновляет меня на то, чтобы продолжать!

Есть несколько главных людей, которые помогли сделать эту книгу реальностью. Я хотел бы поблагодарить моего партнера Зару Каршей за помощь в том, чтобы мой голос был услышан. Спасибо моему редактору Анне Мосс и выпускающему редактору Стефану Лещуку, чьи отзывы и рекомендации были основополагающими в процессе написания, и, в более общем плане, всей редакционной команде Kogan Page's за их усердие и придирчивость. Спасибо моему другу и деловому партнеру Хейделину де Понтевесу за вдохновение и поддержку

в решении некоторых сложных вопросов в области науки о данных, а также за его помощь в рассмотрении технических аспектов этой книги. Мы знаем друг друга чуть больше года, но благодаря нашей дружбе и взаимной поддержке кажется, что знакомы гораздо дольше.

Спасибо Катерине Андрысковой — я чувствую себя обязанным ей и благодарю за то, что она была первым человеком, который прочел «Работу с данными в любой сфере». Моя благодарность талантливой команде SuperDataScience, взявшей на себя множество дополнительных обязанностей, чтобы я мог написать эту книгу. Спасибо трудолюбивой команде Udemy, включая моих неутраченных менеджеров по работе с клиентами Эрин Адамс и Лану Мартинес.

Спасибо моему другу и наставнику Артему Владимирову, чье потрясающее отношение к работе и знания заложили фундамент, на котором я построил все, что знаю сейчас о науке о данных. Огромное спасибо Виталию Долгову и Ричарду Хопкинсу за отличный пример для подражания, за веру в меня, за предоставление помощи всегда, когда я в ней нуждался, за то, что были со мной и в хорошие, и в плохие времена. Спасибо Патрисии Йелленевич, дорогой подруге, преподавшей мне теорию цвета.

Я выражаю благодарность людям, чьи истории использовал, чтобы проиллюстрировать выводы, представленные в этой книге: Хейделину де Понтевесу, Артему Владимирову, Ричарду Хопкинсу, Рубену Коугелу, Грегу Поппу, Кэролайн Макколл, Дэну Шиблеру, Меган Патни и Элфу Морису.

Преподавателей московской школы №54, Московского физико-технического института и Университета Квинсленда я благодарю за то, что получил такое полезное образование. Спасибо всем моим коллегам, бывшим и нынешним, в Deloitte и Sunsuper, за профессиональное развитие, необходимое для создания моего инструментария в науке о данных.

Больше всего я хочу поблагодарить вас, читатель, за то, что уделите мне свое драгоценное время. Я в первую очередь надеюсь, что эта книга поможет тем, кто хочет понять и использовать науку о данных в своей профессиональной деятельности.

Литература

01. Определение данных

Fowles, J. (1992). *Starstruck: Celebrity performers and the American public*, Smithsonian Institution Press.

Keynes, John Maynard (1963). *Essays in Persuasion*, Norton, pp. 358–373.

Mishra, S. and Sharma, M. (2016). Bringing analytics into Indian film industry with back tracing algorithm, *Analytics Vidhya* [Online]: www.analyticsvidhya.com/blog/2016/08/bringing-analytics-into-indian-film-industry-with-back-tracing-algorithm/ [accessed 20.06.2017].

Ohmer, S. (2012). *George Gallup in Hollywood*, Columbia University Press.

02. Как данные удовлетворяют наши потребности

Conley, C. (2007). *Peak: How great companies get their mojo from Maslow*, John Wiley.

Devitt, E. (2017). Artificial chicken grown from cells gets a taste test — but who will regulate it? *Science*, 15 March [Online]: www.sciencemag.org/news/2016/08/lab-grown-meat-inches-closer-us-market-industry-wonders-who-will-regulate [accessed 20.05.2017].

The Economist (2017). The world's most valuable resource is no longer oil, but data, 6 May [Online]: www.economist.com/news/leaders/21721656-data-economy-demands-new-approach-antitrust-rules-worlds-most-valuable-resource [accessed 01.06.2017].

Haselton, T. (2017). Credit reporting firm Equifax says data breach could potentially affect 143 million US consumers. *CNBC*, 7 September. Available from: www.cnn.com/2017/09/07/credit-reporting-firm-equifax-says-cybersecurity-incident-could-potentially-affect-143-million-us-consumers.html [Last accessed 29.09.2017].

IBM (2017a). *Green Horizon: driving sustainable development* [Online]: www.research.ibm.com/labs/china/greenhorizon.html [accessed 12.06.2017].

IBM (2017b). *Watson* [Online]: www.ibm.com/watson/ [accessed 15.02.2017].

Mateos-Garcia, J. and Gardiner, J. (2016). From detecting to engaging: an analysis of emerging tech topics using Meetup data. *Nesta*, 1 August [Online]: www.nesta.org.uk/blog/find-emerging-tech-topics-with-meetup-data [accessed 12.06.2017].

Otake, T. (2016). IBM big data used for rapid diagnosis of rare leukemia case in Japan, *Japan Times*, 11 August [Online]: www.japantimes.co.jp/news/2016/08/11/national/science-health/ibm-big-data-used-for-rapid-diagnosis-of-rare-leukemia-case-in-japan/ [accessed 01.02.2017].

SuperDataScience (2016). SDS 005: Computer forensics, fraud analytics and knowing when to take a break with Dmitry Korneev [Podcast]. 29 March [Online]: www.superdatascience.com/5 [accessed 05.06.17].

03. Мышление, необходимое для эффективного анализа данных

British Academy (2017). Do we need robot law? [Video] [Online]: www.britac.ac.uk/video/do-we-need-robot-law [accessed 03.04.2017].

British Academy and The Royal Society (2017). 'Data management and use: governance in the 21st century' [Online]: www.britac.ac.uk/sites/default/files/Data%20management%20and%20use%20-%20Governance%20in%20the%2021st%20century.pdf [accessed 07.07.2017].

DeepMind (2016). 'DeepMind AI Reduces Google Data Centre Cooling Bill by 40%' [Online]: www.deepmind.com/blog/deepmind-ai-reduces-google-data-centre-cooling-bill-40 [accessed 29.09.2017].

DeepMind (2017). 'AlphaGo Zero: Learning from scratch' [Online]: www.deepmind.com/blog/alphago-zero-learning-scratch [accessed 19.11.2017].

Mulgan, G. (2016). A machine intelligence commission for the UK: how to grow informed public trust and maximise the positive impact of smart machines, *Nesta* [Online]: www.nesta.org.uk/sites/default/files/a_machine_intelligence_commission_for_the_uk_-_geoff_mulgan.pdf [accessed 25.05.17].

Service, R. (2017). DNA could store all of the world's data in one room, *Science*, 2 March [Online]: www.sciencemag.org/news/2017/03/dna-could-store-all-worlds-data-one-room [accessed 27.03.17].

Stallkamp, J. et al. (2012). Man vs computer: benchmark learning algorithms for traffic sign recognition [Online]: <http://image.diku.dk/igel/paper/MvCBMLAFTSR.pdf> [accessed 22.09.17].

UK Cabinet Office (2016). Data science ethical framework [Online]: www.gov.uk/government/publications/data-science-ethical-framework [accessed 07.07.2017].

04. Сформулируйте вопрос

SuperDataScience (2016). SDS 016: Data-driven operations, consulting approaches and mentoring with Richard Hopkins [Podcast]. 22 December [Online]: www.superdatascience.com/16 [accessed 01.08.17].

SuperDataScience (2017). SDS 039: Key data science and statistical skills to get hired at VSCO [Podcast]. 29 March [Online]: www.superdatascience.com/39 [accessed 05.06.17].

05. Подготовка данных

SuperDataScience (2016). SDS 008: data science in computer games, learning to learn and a 40m euro case study with Ulf Morys [Podcast]. 28 October [Online]: www.superdatascience.com/8 [accessed 05.06.17].

06. Анализ данных (часть I)

SuperDataScience (2016a). Data Science A — Z: Real-Life Data Science Exercises Included [Online Course]. Available from: www.udemy.com/datascience/ [Last accessed: 05.06.2017].

SuperDataScience (2016b). SDS 007: Advanced Analytics, Dynamic Simulations, and Consulting Round The Globe with Artem Vladimirov [Podcast]. 21 October. Available from: www.superdatascience.com/39 [Last accessed 05.06.17].

SuperDataScience (2017). Machine Learning A — Z: Hands-On Python & R In Data Science [Online Course]. Available from: www.udemy.com/machinelearning/ [Last accessed: 05.06.2017].

07. Анализ данных (часть II)

Recode (2017). These AI bots created their own language to talk to each other [Online]: www.recode.net/2017/3/23/14962182/ai-learning-language-open-ai-research [accessed 20.11.2017].

08. Визуализация данных

SuperDataScience (2016). SDS 012: Online learning, tableau insights and ad hoc analytics with Megan Putney [Podcast]. 22 November [Online]: www.superdatascience.com/12 [accessed 05.06.2017].

09. Презентация данных

SuperDataScience (2016a). SDS 010: Model validation, data exhaust and organisational cultural change with Yaw Tan [Podcast] 10 November [Online]:

<https://www.superdatascience.com/sds-010-model-validation-data-exhaust-and-organisational-cultural-change-with-yaw-tan> [accessed 05.06.17].

SuperDataScience (2016b). SDS 014: Credit scoring models, the law of large numbers and model building with Greg Poppe [Podcast]. 12 December [Online]: <https://www.superdatascience.com/sds-014-credit-scoring-models-the-law-of-large-numbers-and-model-building-with-greg-poppe> [accessed 05.06.17].

10. Ваша карьера в науке о данных

Каку М. Будущее разума. — М.: Альпина нон-фикшн, 2018.

Baeck, P. (2015). Data for Good: How big and open data can be used for the common good, Nesta [Online]: www.nesta.org.uk/publications/data-good [accessed 05.07.2017].

Burning Glass Technologies and IBM (2017). The Quant Crunch: How the demand for data science skills is disrupting the job market [Online]: www-01.ibm.com/common/ssi/cgi-bin/ssialias?htmlfid=IML14576USEN& [accessed 26.06.17].

Crowdfunder (2017). The Data Scientist Report [Online]: visit.crowdfunder.com/WC-2017-Data-Science-Report_LP.html?src=Website&medium=Carousel&campaign=DSR2017&content=DSR2017 [accessed 01.07.17].

Mannes, J. (2017). Airbnb is running its own internal university to teach data science, Techcrunch, 24 May [Online]: www.techcrunch.com/2017/05/24/airbnb-is-running-its-own-internal-university-to-teach-data-science/ [accessed 26.06.17].

Rosenblum, C. (2017). Hillbillies who code: the former miners out to put Kentucky on the tech map, Guardian, 21 April [Online]: www.theguardian.com/us-news/2017/apr/21/tech-industry-coding-kentucky-hillbillies [accessed 26.06.17].

SuperDataScience (2017a). SDS 059: Changing human behaviour through a driving app [podcast]. 7 June [Online]: www.superdatascience.com/59 [accessed 26.06.17].

SuperDataScience (2017b). SDS 049: Great Tips On Building a Successful Analytics Culture [Podcast]. 13 July. Available from: www.superdatascience.com/49 [Last accessed 10.10.17].

Symons, T. (2016). Councils and the data revolution: 7 ways local authorities can get more value from their data, Nesta, 15 July [Online]: www.nesta.org.uk/blog/councils-and-the-data-revolution-7-ways-local-authorities-can-get-more-value-from-their-data/ [accessed 26.06.17].

Еременко Кирилл

РАБОТА С ДАННЫМИ В ЛЮБОЙ СФЕРЕ

Как выйти на новый уровень,
используя аналитику

Главный редактор *С. Турко*

Руководитель проекта *Л. Разживайкина*

Корректоры *Е. Аксёнова, М. Смирнова*

Компьютерная верстка *М. Поташкин*

Художественное оформление и макет *Ю. Буга*

Дизайн обложки *Ю. Буга*

Подписано в печать 29.05.2019. Формат 70×100/16.

Бумага офсетная № 1. Печать офсетная.

Объем 19 печ.л. Тираж 2000 экз. Заказ № .

ООО «Альпина Паблишер»

123060, Москва, а/я 28

Тел. +7 (495) 980-53-54

www.alpina.ru

e-mail: info@alpina.ru

Знак информационной продукции
(Федеральный закон № 436-ФЗ от 29.12.2010 г.)



Отпечатано в типографии Полиграфическо-издательского комплекса «Идел-Пресс»

филиала АО «ТАТМЕДИА», 420066, г. Казань,

ул. Декабристов, 2. e-mail: id-press@yandex.ru <http://www.idel-press.ru>