

Содержание

ЧАСТЬ IV. НЕОПРЕДЕЛЕННЫЕ ЗНАНИЯ И РАССУЖДЕНИЯ В УСЛОВИЯХ НЕОПРЕДЕЛЕННОСТИ	9
Глава 12. Количественная оценка неопределенности	11
12.1. Действия в условиях неопределенности	11
12.2. Вероятность — основная система обозначений	17
12.3. Логический вывод с использованием полных совместных распределений	27
12.4. Независимость	30
12.5. Правило Байеса и его использование	32
12.6. Наивные байесовские модели	37
12.7. Очередное возвращение в мир вампуса	40
Резюме	45
Библиографические и исторические заметки	46
Упражнения	50
Глава 13. Вероятностные рассуждения	57
13.1. Представление знаний в неопределенной проблемной области	57
13.2. Семантика байесовских сетей	61
13.3. Точный вывод в байесовских сетях	80
13.4. Приближенный вероятностный вывод в байесовских сетях	92
13.5. Причинно-следственные байесовские сети	113
Резюме	119
Библиографические и исторические заметки	120
Упражнения	129
Глава 14. Вероятностные рассуждения во времени	139
14.1. Время и неопределенность	140
14.2. Вероятностный вывод во временных моделях	146
14.3. Скрытые марковские модели	158
14.4. Фильтры Калмана	166

14.5. Динамические байесовские сети	176
Резюме	193
Библиографические и исторические заметки	194
Упражнения	198
Глава 15. Вероятностное программирование	205
15.1. Реляционные вероятностные модели	206
15.2. Вероятностные модели с открытой вселенной	216
15.3. Отслеживание состояния сложного мира	227
15.4. Программы как вероятностные модели	233
Резюме	240
Библиографические и исторические заметки	241
Глава 16. Принятие простых решений	247
16.1. Сочетание убеждений и желаний в условиях неопределенности	248
16.2. Основы теории полезности	249
16.3. Функции полезности	254
16.4. Многоатрибутные функции полезности	265
16.5. Сети принятия решений	272
16.6. Стоимость информации	276
16.7. Неизвестные предпочтения	285
Резюме	291
Библиографические и исторические заметки	292
Упражнения	297
Глава 17. Принятие сложных решений	307
17.1. Задачи последовательного принятия решений	307
17.2. Алгоритмы для задач MDP	323
17.3. Задачи о бандитах	335
17.4. Марковские процессы принятия решений в частично наблюдаемых средах	346
17.5. Алгоритмы для решения задач POMDP	350
Резюме	357
Библиографические и исторические заметки	358
Упражнения	362

Глава. 18. Принятие решений при наличии нескольких агентов	367
18.1. Свойства мультиагентной среды	367
18.2. Теория некооперативных игр	376
18.3. Теория кооперативных игр	407
18.4. Принятие коллективных решений	417
Резюме	437
Библиографические и исторические заметки	438
Упражнения	445
Приложение А. Математические основы	447
А.1. Анализ сложности и нотация $O()$	447
А.2. Векторы, матрицы и линейная алгебра	451
А.3. Распределения вероятностей	453
Библиографические и исторические заметки	457
Приложение Б. Сведения о языках и алгоритмах, используемых в книге	458
Б.1. Определение языков с помощью формы Бэкуса–Наура	458
Б.2. Описание алгоритмов с помощью псевдокода	459
Б.3. Дополнительный материал в Интернете	461
Предметный указатель	463

Количественная оценка неопределенности

В данной главе показано, как должен действовать агент в условиях неопределенности со степенью доверия, представленной в числовом виде.

12.1. Действия в условиях неопределенности

Агенты в реальном мире должны справляться с ► **неопределенностью**, будь то по причине частичной наблюдаемости, недетерминизма или действий противников. Агент может никогда не знать наверняка, в каком состоянии он сейчас находится или где он окажется после выполнения некоторой последовательности действий.

Мы уже знакомы с агентами, решающими задачи, и логическими агентами, справляющимися с неопределенностью посредством отслеживания цепочки **доверительных состояний** — представления множества всех возможных состояний мира, которые могут в нем проявиться — и формирования условного плана, позволяющего справиться с любой возможной случайной ситуацией, о которой его датчики могут сообщить во время его выполнения. Подобный подход работает для простых задач, но имеет определенные недостатки.

- Агент должен рассмотреть *все возможные* объяснения для результатов наблюдений своих датчиков, и при этом не важно, насколько эти объяснения будут маловероятны. Такой подход приводит к формированию огромного доверительного состояния, полного почти невероятных возможностей.
- Правильный условный план, учитывающий любую возможность, может вырастать до сколь угодно больших размеров и должен учитывать сколь угодно маловероятные непредвиденные обстоятельства.
- Иногда не существует плана, который гарантированно приводит к цели, но агент все равно должен действовать. Он должен иметь какой-то способ сравнения достоинств различных планов, которые не являются гарантированно достигающими цели.

Например, предположим, что автоматизированному такси поставлена цель — доставить пассажира в аэропорт к заданному времени. Агент такси формирует план, A_{90} , предусматривающий выезд из дома за 90 минут до установленного времени отправления рейса и движение такси с разумной скоростью. Даже если аэропорт находится всего в 5 милях от дома, логический агент не сможет с абсолютной уверенностью прийти к заключению, что “План A_{90} позволяет добраться до аэропорта к назначенному времени”. Вместо этого он придет к более слабому заключению: “План A_{90} позволяет прибыть в аэропорт вовремя, если машина не сломается и не попадет в аварию, и дорога не будет закрыта, и в машину не попадает метеорит, и...” Ни по одному из этих условий нельзя вынести гарантированно верное суждение, поэтому невозможно сделать вывод, что план обязательно будет успешным. Это логическая **проблема квалификации** (см. раздел 7.7.1), для которой до сих пор не найдено реального решения.

Тем не менее в некотором смысле план A_{90} действительно представляет собой правильное руководство к действию. Что имеется в виду под этим утверждением? Как уже говорилось в главе 2, под этим подразумевается, что из всех планов, которые могут быть выполнены, именно план A_{90} , как ожидается, позволит максимизировать показатели производительности агента (здесь это ожидание строится на основании знаний агента об окружающей среде). Показатели производительности включают своевременную доставку пассажира в аэропорт к указанному рейсу, предотвращение продолжительного, непродуктивного ожидания в аэропорту и исключение штрафов за превышение скорости по пути в аэропорт. Знания агента не позволяют гарантировать достижения любого из этих трех результатов при выполнении плана A_{90} , но могут обеспечить некоторую степень уверенности, что они будут достигнуты. Другие планы, например A_{180} , могут повысить степень уверенности агента в том, что он доставит пассажира до аэропорта вовремя, но одновременно повысят для него и вероятность продолжительного, скучного ожидания.

➔ *Следовательно, выбор правильного способа действия — рационального решения — зависит как от относительной важности различных целей, так и от степени уверенности в том, что они могут быть достигнуты.* В оставшейся части данного раздела эти идеи будут уточнены с целью подготовки к разработке общих теорий проведения рассуждений в условиях неопределенности и принятия рациональных решений, которые будут представлены в этой и последующих главах.

12.1.1. Учет наличия неопределенности

Рассмотрим простой пример рассуждений при наличии неопределенности: диагностика причин зубной боли у пациента. Диагностика — при медицинском обследовании, при ремонте автомобиля или в любых других случаях — почти всегда связана с неопределенностью. Попробуем записать правила для диагностики заболеваний зубов с использованием логики высказываний, что явным образом укажет на трудности, возникающие при простом логическом

подходе. Рассмотрим следующее простое правило (здесь *Toothache* — зубная боль, а *Cavity* — полость):

$$\textit{Toothache} \Rightarrow \textit{Cavity}.$$

Проблема состоит в том, что это правило неверно. Не у всех пациентов с зубной болью обязательно имеется полость, — у некоторых причиной боли может быть заболевание десен (*GumProblem*), нарыв (*Abscess*) или одна из нескольких иных сложных ситуаций.

$$\textit{Toothache} \Rightarrow \textit{Cavity} \vee \textit{GumProblem} \vee \textit{Abscess} \dots$$

К сожалению, чтобы сделать это правило истинным, потребуется ввести в него почти бесконечный список возможных причин. Это правило можно попытаться преобразовать в причинное правило:

$$\textit{Cavity} \Rightarrow \textit{Toothache}.$$

Но и это правило нельзя назвать верным; не все зубы, имеющие полость, обязательно вызывают болевые ощущения. Единственный способ исправить данное правило состоит в том, чтобы сделать его логически исчерпывающим: дополнить левую сторону описаниями всех обстоятельств, которые должны иметь место для того, чтобы полость действительно вызывала зубную боль. Следовательно, попытка использовать логику для решения задач в проблемной области, подобной медицинской диагностике, оканчивается неудачей по следующим трем основным причинам.

- ► **Экономия усилий.** Для формирования полного множества antecedентов или консеквентов, необходимого для составления правила, не имеющего исключений, потребуется слишком много работы, а применение таких правил будет слишком сложным.
- ► **Неполнота теоретических знаний.** Медицинская наука не имеет полной теории для данной проблемной области.
- ► **Неполнота практических знаний.** Даже если известны все теоретические правила, может иметь место неопределенность в отношении диагноза для конкретного пациента, поскольку не все необходимые обследования были или вообще могут быть выполнены.

Связь между зубной болью и наличием полости не является простым логическим следствием, действующим в обоих направлениях. Такая ситуация типична не только для медицинской диагностики, но и для большинства других проблемных областей, связанных с вынесением суждений: юриспруденции, бизнеса и экономики, проектирования, ремонта автомобилей, садоводства, датирования объектов или событий и т.д. Знания агента в лучшем случае позволяют сформулировать релевантные суждения только с определенной ► **степенью уверенности** (*degree of belief*). Нашим основным инструментальным средством для работы со степенью уверенности будет ► **теория вероятностей**. В соответствии с терминологией,

представленной в разделе 8.1, **онтологический вклад** логики и теории вероятностей одинаков — мир состоит из фактов, которые имеют либо не имеют места в каждом конкретном случае, но их **эпистемологический вклад** будет разным: логический агент уверен, что каждое высказывание должно быть истинным или ложным, либо у него нет никакого мнения, в то время как вероятностный агент может иметь числовую оценку степени уверенности в диапазоне от 0 (для высказываний, которые точно ложны) до 1 (для высказываний, которые безусловно верны).

→ *Теория вероятности предоставляет способ суммарного учета неопределенности, возникающей по причинам экономии усилий и неполноты знаний*, тем самым решая проблему квалификации. Можно не знать со всей уверенностью, что именно вызывает зубную боль у определенного пациента, но можно с уверенностью полагать, что, скажем, в 80% случаев — т.е. с вероятностью 0,8, — если пациент испытывает зубную боль, то ее источником является полость в зубе. Это означает, что из всех ситуаций, неотличимых от текущей в пределах тех знаний, которыми обладает агент, в 80% этих случаев у пациента должна быть полость в зубе. Подобная уверенность может быть основана на статистических данных — у 80% пациентов с зубной болью, наблюдавшихся до сих пор, была обнаружена зубная полость, — на некоторых общих знаниях из области стоматологии или на комбинации различных источников.

Один вносящий путаницу момент состоит в том, что при постановке диагноза в реальном мире нет никакой неопределенности: в зубе пациента либо есть полость, либо нет. Так что же означает наше утверждение, что вероятность наличия полости равна 0,8? Разве она не должна быть равна 0 или 1? Ответ состоит в том, что вероятностные высказывания делаются в отношении *состояния знаний* агента, а не в отношении *реального мира*. Мы говорим: “Вероятность того, что пациент имеет зубную полость, *принимая во внимание то, что он испытывает зубную боль*, равна 0,8”. Если позднее выяснится, что пациент уже некоторое время страдает заболеванием десен, можно будет прийти к другому заключению: “Вероятность того, что пациент имеет зубную полость, *принимая во внимание, что он испытывает зубную боль и страдает заболеванием десен*, составляет 0,4”. Если будут собраны дополнительные убедительные доказательства против наличия зубной полости, появится возможность утверждать: “Вероятность того, что пациент имеет зубную полость, с учетом всего того, что нам теперь известно, почти нулевая”. Обратите внимание, что все приведенные выше заключения не противоречат друг другу, — в каждом есть собственное утверждение о различном состоянии знаний агента.

12.1.2. Неопределенность и рациональные решения

Еще раз вернемся к плану поездки в аэропорт A_{90} . Предположим, что он обеспечивает 97%-ный шанс успешного вылета назначенным рейсом. Означает ли это, что выбор данного плана будет рациональным решением? Вовсе необязательно: могут существовать другие планы — например, A_{180} — с большей вероятностью

успешного вылета. Если *жизненно* важно не пропустить именно этот рейс, то стоит рискнуть подождать в аэропорту подольше. А что можно сказать о плане A_{1440} , предусматривающем заблаговременный выезд из дома за 24 часа до отправления самолета? В большинстве ситуаций это будет не лучший выбор, поскольку, хотя он практически гарантирует прибытие в аэропорт вовремя, он предполагает и невыносимо долгое ожидание, не говоря уже о возможности малоприятной диеты из предлагаемого в аэропорту меню.

Чтобы сделать подобный выбор, агент прежде всего должен иметь сведения о **▶ предпочтениях** среди различных возможных **▶ результатов** различных планов. Любой результат — это полностью определенное состояние, включая такие факторы, как своевременность прибытия и длительность ожидания в аэропорту. Для представления предпочтений и количественных рассуждений о них мы будем использовать **▶ теорию полезности** (*utility theory*). (Термин “полезность” в данном контексте обозначает “свойство быть полезным”.) Теория полезности утверждает, что для агента каждое состояние (или последовательность состояний) имеет определенную степень полезности (или просто *полезность*) и что агент всегда должен отдавать предпочтение состояниям с более высокой полезностью.

Для агента полезность состояния является величиной относительной. Например, полезность состояния, в котором белые могут поставить мат черным при игре в шахматы, очевидно высока для агента, играющего белыми, и очень низка для агента, играющего черными. Но мы не можем строго следовать оценкам в 1, 1/2 и 0 баллов, которые диктуются правилами проведения шахматных турниров, — одни игроки (включая авторов книги) могут быть в восторге от ничьей с чемпионом мира, тогда как другие игроки (включая прежнего чемпиона мира), едва ли будут ей особенно рады. В любом случае личные вкусы или предпочтения не должны учитываться: можно полагать, что агент, отдающий предпочтение мороженому с вкраплениями жевательной резинки “Халапеньо” вместо изюма или шоколадных чипсов, — очень странный, но нельзя утверждать, что он очевидно нерационален. Функция полезности может учитывать любое множество предпочтений — необычных или типичных, благородных или порочных. Можно даже учитывать полезность альтруистического поведения, просто включив оценку благополучия других как один из факторов.

Предпочтения, выраженные в виде полезности, комбинируются с вероятностями в общей теории рациональных решений, называемой **▶ теорией принятия решений** (*decision theory*):

Теория принятия решений = теория вероятностей + теория полезности.

Фундаментальная идея теории принятия решений состоит в том, что **➔** *любой агент является рациональным тогда и только тогда, когда он выбирает действие, позволяющее достичь наибольшей ожидаемой полезности, усредненной по всем возможным результатам данного действия.* Это — принцип **▶ максимальной ожидаемой полезности** (*Maximum Expected Utility* — MEU). Здесь “ожидаемой” означает

“средней”; точнее, это “статистическое среднее” значений полезностей, взвешенных по вероятности их получения. Мы наблюдали этот принцип в действии в главе 5, когда обсуждали выбор оптимальных решений при игре в нарды. В действительности это совершенно общий принцип принятия решений для агентов, действующих в одиночку.

На рис. 12.1 приведен набросок структуры агента, использующего теорию принятия решений для выбора действия. На некотором абстрактном уровне этот агент идентичен агентам, описанным в главах 4 и 7, которые поддерживают доверительное состояние, отражающее историю восприятий на текущий момент. Основное различие заключается в том, что доверительное состояние агента, действующего в соответствии с теорией принятия решений, представляет не только *возможности* для состояний мира, но и их *вероятности*. Основываясь на доверительном состоянии и некоторых знаниях о результатах действий, агент может сделать вероятностные предсказания о результатах выполнения действия и, следовательно, выбрать действие с наибольшей ожидаемой полезностью.

function DT-AGENT(*percept*) **returns** действие *action*
persistent: *belief_state*, доверительное состояние — вероятностные
убеждения в отношении текущего состояния мира
action, действие агента

обновить *belief_state* с учетом действия *action* и восприятия *percept*
вычислить результирующие вероятности для действий *actions*
на основании описаний действий *action* и текущего доверительного
состояния *belief_state*
выбрать действие *action* с наивысшей ожидаемой полезностью,
исходя из вероятностей результатов и информации о полезности
return *action*

Рис. 12.1. Агент, действующий на основании теории принятия решений и выбирающий рациональные действия

В этой и следующей главах изложение в основном сосредоточено на задаче представления данных и вычислений с учетом вероятностной информации в целом. Глава 14 посвящена методам решения конкретных задач представления и обновления доверительного состояния во времени и прогнозированию результатов. В главе 15 рассматриваются способы комбинирования теории вероятностей с выразительными формальными языками, такими как логика первого порядка и языки программирования общего назначения. В главе 16 теория полезности рассматривается более подробно, а в главе 17 разрабатываются алгоритмы планирования последовательностей действий в стохастических средах. В главе 18 все эти идеи распространяются на многоагентные проблемные среды.

12.2. Вероятность — основная система обозначений

Чтобы агент мог представлять и использовать вероятностную информацию, необходим формальный язык представления неопределенных знаний. Язык теории вероятностей традиционно является неформальным, разработанным человеком-математиком для других математиков. Стандартное введение в элементарную теорию вероятностей вы найдете в приложении А. В этом разделе выбран иной подход, более удобный для потребностей ИИ, в котором теория вероятностей соединяется с понятиями формальной логики.

12.2.1. О каких вероятностях идет речь

Подобно логическим утверждениям, вероятностные утверждения относятся к возможным мирам. В то время как логические утверждения говорят, какие из возможных миров являются строго недопустимыми (все те, в которых утверждение является ложным), вероятностные утверждения говорят о том, насколько вероятными являются различные миры. В теории вероятностей множество всех возможных миров называют ► **пространством элементарных событий**. Возможные миры являются *взаимоисключающими* и *исчерпывающими* — два возможных мира не могут иметь место одновременно, иметь место в реальности допустимо лишь для одного возможного мира. Например, если мы собираемся бросить две (отличимые друг от друга) кости, существует 36 возможных миров, которые следует рассмотреть: (1,1), (1,2), ..., (6,6). Для обозначения пространства элементарных событий используется греческая буква Ω (прописная буква “омега”), а буква ω (строчная буква “омега”) используется для ссылок на элементы этого пространства, т.е. на конкретные возможные миры, которые в данном контексте также называют ► **элементарными событиями**.

Полностью определенная ► **вероятностная модель** связывает числовую вероятность $P(\omega)$ с каждым возможным миром.¹ Основные аксиомы теории вероятностей говорят о том, что каждый возможный мир характеризуется вероятностью в пределах от 0 до 1 и что суммарная вероятность всего множества возможных миров равна 1.

$$0 \leq P(\omega) \leq 1 \quad \text{для каждого } \omega \text{ и } \sum_{\omega \in \Omega} P(\omega) = 1 \quad (12.1)$$

Например, если предположить, что каждая кость выполнена без изъянов и при броске они не мешают друг другу, то каждый из возможных миров (1,1), (1,2), ..., (6,6) характеризуется вероятностью 1/36. Если некоторые грани костей будут

¹ На данный момент мы предполагаем множество возможных миров дискретным и счетным. Корректная обработка непрерывных множеств требует учета определенных сложных моментов, которые менее актуальны для большинства целей в ИИ.

дополнительно утяжелены, то одни миры будут иметь более высокую вероятность, а другие — более низкую, но общая их сумма все равно составит 1.

Вероятностные утверждения и запросы обычно касаются не конкретных возможных миров, а некоторых их множеств. Например, нас может интересовать вероятность того, что при броске двух костей сумма будет равна 11, или вероятность того, что выпадут два одинаковых значения, и т.д. В теории вероятностей эти множества называются ► **событиями** или **исходами** — этот термин уже широко использовался в главе 10 для иной концепции. В логике множество миров соответствует **высказыванию** на формальном языке, в частности для каждого высказывания соответствующее множество содержит только те возможные миры, в которых это высказывание истинно. (Таким образом, в этом контексте “событие” и “высказывание” означают примерно одно и то же, за исключением того, что высказывание выражается средствами формального языка.) Вероятность, связанная с высказыванием, определяется как сумма вероятностей тех возможных миров, в которых оно истинно.

$$\text{Для любого высказывания } \varphi, P(\varphi) = \sum_{\omega \in \Omega} P(\omega). \quad (12.2)$$

Например, при броске симметричных костей вероятность выпадения суммы 11 можно определить как $P(\text{Total} = 11) = P((5,6)) + P((6,5)) = 1/36 + 1/36 = 1/18$. Следует отметить, что теория вероятности не требует полного знания вероятностей для каждого возможного мира. Например, если мы полагаем, что кости были переделаны так, чтобы при броске на них чаще выпадали одинаковые значения, то можем *утверждать*, что $P(\text{doubles}) = 1/4$, даже не принимая во внимание то, что при броске чаще будут выпадать, скажем, две шестерки, а не две двойки. Так же, как и в случае логических утверждений, это утверждение *ограничивает* лежащую в основе вероятностную модель без полного ее определения.

Вероятности, такие как $P(\text{Total} = 11)$ и $P(\text{doubles})$, принято называть ► **безусловными** или ► **априорными вероятностями**, — они касаются степени доверия к высказываниям *при отсутствии какой-либо другой информации*. Чаще всего, однако, у нас есть *некоторая* информация, обычно называемая ► **свидетельством**, которая уже была получена ранее. Например, при броске двух костей первая из них уже может остановиться со значением 5 и мы, затаив дыхание, ждем, когда остановится вторая. В этом случае нас интересует не априорная вероятность результата броска костей, а ► **условная** или ► **апостериорная вероятность** выпадения дубля, *когда на первой кости (Die₁) уже выпала пятерка*. Эта вероятность записывается как $P(\text{doubles} | \text{Die}_1 = 5)$, где символ “|” читается как “при условии”.²

Аналогичным образом, если отправиться к стоматологу на регулярное плановое обследование, то априорная вероятность $P(\text{cavity}) = 0,2$ может представлять

² Оператор | имеет наименьший возможный приоритет, поэтому выражение $P(\dots | \dots)$ всегда означает $P(\dots)(\dots)$.

интерес, но если отправиться стоматологу потому, что болит зуб, то важнее будет апостериорная вероятность $P(\text{cavity} | \text{toothache}) = 0,6$.

Важно понимать, что вероятность $P(\text{cavity}) = 0,2$ по-прежнему остается *имеющей силу* и после появления *зубной боли*, просто она уже не является особенно полезной. Принимая решения, агенту нужно обусловить *все* те свидетельства, которые он наблюдал. Также важно понимать различие, существующее между обусловливанием и логическим следствием. Утверждение, что $P(\text{cavity} | \text{toothache}) = 0,6$, не означает “Всякий раз, когда имеет место *зубная боль*, можно сделать заключение, что наличие в зубе *полости* имеет место с вероятностью 0,6”. В действительности оно означает “Всякий раз, когда имеет место *зубная боль* и у нас нет никакой *дополнительной информации*, можно сделать заключение, что наличие в зубе *полости* имеет место с вероятностью 0,6”. Это дополнительное обусловливание имеет большое значение. Например, получив дополнительную информацию о том, что стоматолог так и не обнаружил в зубах пациента никаких полостей, мы, определенно, не захотим прийти к заключению, что наличие в зубе полости имеет вероятность 0,6; вместо этого нам нужно будет использовать заключение $P(\text{cavity} | \text{toothache} \wedge \neg \text{cavity}) = 0$.

Говоря языком математики, апостериорные вероятности определяются в терминах априорных вероятностей следующим образом. Для любых высказываний a и b мы имеем

$$P(a | b) = \frac{P(a \wedge b)}{P(b)}, \quad (12.3)$$

что выполняется всякий раз, когда $P(b) > 0$. Например,

$$P(\text{doubles} | \text{Die}_1 = 5) = \frac{P(\text{doubles} \wedge \text{Die}_1 = 5)}{P(\text{Die}_1 = 5)}.$$

Это определение обретает смысл, если вспомнить, что наблюдение правил b исключает все те возможные миры, в которых b является ложным, оставляя множество, суммарная вероятность которого — просто $P(b)$. В пределах этого множества миры, в которых a является истинным, должны удовлетворять условию $a \wedge b$ и представляют собой дробь $P(a \wedge b)/P(b)$.

Определение апостериорной (или условной) вероятности — уравнение (12.3) — может быть записано в другой форме, называемой **правилем умножения вероятностей**:

$$P(a \wedge b) = P(a | b)P(b). \quad (12.4)$$

Правило умножения вероятностей, возможно, легче запомнить: оно построено на основании того факта, что для того, чтобы $a \wedge b$ было истинно, необходимо, чтобы b было истинно, а также чтобы истинно было a при данном b .

12.2.2. Язык высказываний в вероятностных утверждениях

В этой и последующих главах высказывания, описывающие множества возможных миров, как правило, будут записываться с использованием нотации, в которой сочетаются элементы логики высказываний, и нотации удовлетворения ограничений. Согласно терминологии, предложенной в разделе 2.4.7, это **развернутое представление**, в котором возможный мир представлен в виде множества пар *переменная/значение*. Также возможно использование еще более выразительного **структурного представления**, как показано в главе 15.

В теории вероятностей переменные называются ► **случайными переменными** и их имена всегда начинаются с прописной буквы. Таким образом, в примере с бросанием костей переменные *Total* и *Die₁* являются случайными. Каждая случайная переменная является функцией, отображающей проблемную область возможных миров Ω в некоторую ► **область определения значений** (*range*) — множество возможных значений, которые она может принимать. Областью определения переменной *Total* в случае двух костей является множество $\{2, \dots, 12\}$, а областью определения переменной *Die₁* — $\{1, \dots, 6\}$. Имена значений всегда записываются строчными буквами, поэтому сумму всех значений переменной X можно записать как $\sum_x P(X=x)$. Булева случайная переменная имеет область определения $\{true, false\}$. Например, высказывание о том, что были выброшены дубли, можно записать как *Doubles = true*. (Альтернативным вариантом области определения для булевых переменных является множество $\{0, 1\}$, и в этом случае говорят, что переменная имеет распределение ► **Бернулли**.) По соглашению высказывания вида $A = true$ сокращаются до просто a , тогда как высказывания вида $A = false$ сокращаются до $\neg a$. (Использованные в предыдущем разделе имена *doubles*, *cavity* и *toothache* являются сокращением именно этого типа.)

Области определения значений переменных могут представлять собой множество произвольных признаков. Например, для переменной *Age* (*возраст*) можно выбрать область определения в виде множества $\{juvenile, teen, adult\}$ (т.е. *ребенок, подросток, взрослый*), а для переменной *Weather* (*погода*) областью определения могут быть значения $\{sun, rain, cloud, snow\}$ (т.е. *солнечно, дождь, облачно, снег*). Если неоднозначное понимание исключено, то обычно принято использовать само значение в тех высказываниях, где определенная переменная имеет это значение; так, значение *sun* можно непосредственно использовать в высказывании $Weather = sun$.³

Все предыдущие примеры имеют конечные области определения значений. Переменные также могут иметь бесконечные области определения — либо

³ Эти соглашения, взятые вместе, приводят к потенциальной неоднозначности в обозначениях при суммировании значений булевых переменных. Например, $P(a)$ — вероятность того, что переменная A имеет значение *true*, тогда как в выражении $\sum_a P(a)$ это просто ссылка на вероятность одного из значений a переменной A .

дискретные (например, целые числа), либо непрерывные (например, действительные числа). Для любой переменной с упорядоченной областью определения также допускаются неравенства, такие как $NumberOfAtomsInUniverse \geq 10^{70}$.

Наконец, можно объединить все эти виды элементарных высказываний (включая сокращенные формы для булевых переменных), используя стандартные логические связки логики высказываний. Например, высказывание “Вероятность того, что в зубах пациентки есть полость, с учетом того, что она является подростком и не испытывает зубной боли, составляет 0,1” можно записать следующим образом.

$$P(cavity | \neg toothache \wedge teen) = 0,1$$

Также в вероятностной нотации для обозначения операции конъюнкции часто используют запятую, поэтому в приведенном выше высказывании левую часть можно было бы записать просто как $P(cavity | \neg toothache, teen)$.

Иногда в обсуждение требуется включить вероятности *всех* возможных значений случайной величины. Понятно, что в этом случае можно было бы использовать такую запись:

$$\begin{aligned} P(Weather = sun) &= 0,6, \\ P(Weather = rain) &= 0,1, \\ P(Weather = cloud) &= 0,29, \\ P(Weather = snow) &= 0,01, \end{aligned}$$

но для сокращения можно применить следующий вариант записи:

$$P(Weather) = (0,6; 0,1; 0,29; 0,01).$$

Здесь выделение **P** полужирным шрифтом указывает, что результатом является вектор чисел, расположенных в некотором предопределенном порядке $\langle sun, rain, cloud, snow \rangle$ в соответствии с областью определения переменной *Weather*. Говорят, что высказывание **P** задает ► **распределение вероятностей** для случайной переменной *Weather*, т.е. присвоение вероятности для каждого возможного значения этой случайной переменной. (В подобном случае при конечной дискретной области определения значений такое распределение называется ► **категориальным распределением**.) Нотация **P** также используется для условных распределений: $P(X|Y)$, присваивая значения $P(X=x_i | Y=y_j)$ для каждой возможной пары i, j .

Для непрерывных переменных просто невозможно записать все распределение в виде вектора, поскольку в нем существует бесконечно много значений. Вместо этого можно определить вероятность того, что случайная величина принимает некоторое значение x как параметризованную функцию от x , обычно называемую ► **функцией плотности распределения вероятностей**. Например, высказывание

$$P(NoonTemp = x) = Uniform(x; 18C, 26C)$$

выражает уверенность в том, что значение температуры в полдень (переменная $NoonTemp$) будет равномерно распределено между значениями 18 и 26 градусов по Цельсию.

Функция плотности распределения вероятностей (также часто называемая просто **функцией распределения вероятностей**) по смыслу отличается от дискретных распределений. Утверждение, что плотность вероятности равномерно распределена от 18°C до 26°C, означает, что существует 100%-ная вероятность того, что значение температуры в полдень попадет в этот диапазон шириной в 8°C, и 50%-ная вероятность того, что оно попадет в любой поддиапазон шириной 4°C этого диапазона, и т.д. Принято записывать плотность вероятности для непрерывной случайной величины X в области значения x как $P(X=x)$ или просто как $P(x)$. Интуитивно понятное определение $P(x)$ — это вероятность того, что значение X попадает в произвольно малую область, начинающуюся от x , деленную на ширину этой области:

$$P(x) = \lim_{dx \rightarrow 0} P(x \leq X \leq x + dx) / dx.$$

Для переменной $NoonTemp$ имеем

$$P(NoonTemp = x) = Uniform(x; 18C, 26C) = \begin{cases} \frac{1}{8C} & \text{если } 18C \leq x \leq 26C \\ 0 & \text{в противном случае.} \end{cases}$$

Здесь C обозначает шкалу температуры в градусах Цельсия (а не является константой). В выражении $P(NoonTemp = 20, 18C) = \frac{1}{8C}$ обратите внимание, что $\frac{1}{8C}$ — это не вероятность, а *плотность вероятности*. Вероятность того, что переменная $NoonTemp$ имеет значение *точно* 20,18°C, равна нулю, потому что диапазон 20,18°C имеет нулевую ширину. Некоторые авторы используют разные символы для дискретных вероятностей и плотностей вероятностей; но мы в этой книге будем использовать обозначение P для конкретных значений вероятности и \mathbf{P} — для векторов значений в обоих случаях, поскольку в действительности путаница возникает очень редко и уравнения чаще всего идентичны. Обратите внимание, что вероятности — это безразмерные числа, тогда как значения функций распределения вероятностей выражаются в некоторых единицах измерения. В нашем примере это единица, обратная градусу Цельсия. Если тот же самый интервал температур потребуется выразить в градусах по Фаренгейту, он будет иметь ширину 14,4 градуса, а плотность вероятности будет иметь значение $1/14,4F$.

В дополнение к распределениям по отдельным переменным нам потребуются обозначения и для распределений по нескольким переменным. Для этой цели будем использовать запятые. Например, выражение $\mathbf{P}(Weather, Cavity)$ определяет вероятности всех комбинаций значений переменных $Weather$ и $Cavity$ и представляет собой таблицу вероятностей размером 4×2 , называемую **совместным рас-**

пределием вероятностей для переменных *Weather* и *Cavity*. Также можно смешивать в выражениях переменные и конкретные значения, например $\mathbf{P}(sun, Cavity)$ является двухэлементным вектором, включающим вероятности наличия в зубе полости в солнечный день и отсутствия полости в зубе в солнечный день.

Использование обозначения \mathbf{P} делает определение выражений гораздо более кратким, чем они могли бы быть в противном случае. Например, правило умножения вероятностей (см. уравнение (12.4)) для всех возможных значений переменных *Weather* и *Cavity* можно записать в виде единственного уравнения:

$$\mathbf{P}(Weather, Cavity) = \mathbf{P}(Weather | Cavity)\mathbf{P}(Cavity)$$

вместо следующих $4 \times 2 = 8$ уравнений (с использованием сокращений *W* и *C*):

$$\begin{aligned} P(W = sun \wedge C = true) &= P(W = sun | C = true) P(C = true) \\ P(W = rain \wedge C = true) &= P(W = rain | C = true) P(C = true) \\ P(W = cloud \wedge C = true) &= P(W = cloud | C = true) P(C = true) \\ P(W = snow \wedge C = true) &= P(W = snow | C = true) P(C = true) \\ P(W = sun \wedge C = false) &= P(W = sun | C = false) P(C = false) \\ P(W = rain \wedge C = false) &= P(W = rain | C = false) P(C = false) \\ P(W = cloud \wedge C = false) &= P(W = cloud | C = false) P(C = false) \\ P(W = snow \wedge C = false) &= P(W = snow | C = false) P(C = false). \end{aligned}$$

Как вырожденный случай выражение $\mathbf{P}(sun, cavity)$ не содержит переменных и, следовательно, является вектором нулевой размерности, который можно рассматривать как скалярное значение.

На данный момент мы уже определили синтаксис высказываний и вероятностных утверждений, а также дали часть семантики: уравнение (12.2) определяет вероятность высказывания как сумму вероятностей миров, в которых оно выполняется. Для завершения семантики необходимо сказать, чем эти миры являются и как определить, выполняется ли высказывание в некотором мире. Мы заимствуем эту часть непосредственно из семантики логики высказываний следующим образом. ➔ *Возможный мир определяется как присваивание значений всем рассматриваемым случайным переменным.*

Легко показать, что это определение удовлетворяет основному требованию, согласно которому возможные миры должны быть взаимно исключающими и исчерпывающими (см. упражнение 12.5). Например, если случайными переменными являются *Cavity*, *Toothache* и *Weather*, то существует $2 \times 2 \times 4 = 16$ возможных миров. Более того, истинность любого заданного высказывания легко может быть определена в подобных мирах посредством тех же самых рекурсивных вычислений истинности, которые использовались нами в логике высказываний (см. раздел 7.4.2).

Обратите внимание, что некоторые случайные переменные могут быть избыточными в том смысле, что их значения во всех случаях могут быть получены из значений других переменных. Например, в мире двух игральные

костей переменная *Doubles* будет иметь значение *true* только в тех случаях, когда $Die_1 = Die_2$. Включение переменной *Doubles* в качестве одной из случайных переменных в дополнение к переменным Die_1 и Die_2 , как кажется, увеличивает количество возможных миров с 36 до 72, но, конечно же, ровно половина из этих 72 миров будет логически невозможной и, следовательно, иметь вероятность 0.

Из приведенного выше определения возможных миров следует, что вероятностная модель полностью определяется совместным распределением вероятностей для всех случайных переменных — так называемым ► **полным совместным распределением вероятностей**. Например, при наличии случайных переменных *Cavity*, *Toothache* и *Weather* полным совместным распределением вероятностей будет $P(Cavity, Toothache, Weather)$. Это совместное распределение может быть представлено в виде таблицы размерностью $2 \times 2 \times 4$, содержащей 16 значений. Поскольку вероятность каждого высказывания является суммой по всем возможным мирам, полного совместного распределения в принципе достаточно для вычисления вероятности любого высказывания. Примеры того, как это можно сделать, будут приведены в разделе 12.3.

12.2.3. Аксиомы вероятности и их обоснованность

Основные аксиомы вероятности (уравнения (12,1) и (12,2)) подразумевают определенные отношения между степенями доверия, которые могут быть отнесены к логически связанным высказываниям. Так, можно вывести знакомые отношения между вероятностью высказывания и вероятностью его отрицания:

$$\begin{aligned}
 P(\neg a) &= \sum_{\omega \in \neg a} P(\omega) = && \text{по уравнению (12.2)} \\
 &= \sum_{\omega \in \neg a} P(\omega) + \sum_{\omega \in a} P(\omega) - \sum_{\omega \in a} P(\omega) = \\
 &= \sum_{\omega \in \Omega} P(\omega) - \sum_{\omega \in a} P(\omega) = && \text{группируя первые 2 члена} \\
 &= 1 - P(a) && \text{по (12.1) и (12.2).}
 \end{aligned}$$

Также можно вывести известную формулу для вероятности дизъюнкции, которую иногда называют ► **формулой (или принципом) включений-исключений**:

$$P(a \vee b) = P(a) + P(b) - P(a \wedge b). \tag{12.5}$$

Это правило можно легко запомнить, отметив, что те случаи, когда высказывание *a* является истинным, вместе с теми случаями, когда высказывание *b* является истинным, безусловно, охватывают все те случаи, когда истинно высказывание $a \vee b$; но в сумме двух множеств случаи их пересечения будут учтены дважды, поэтому необходимо вычесть $P(a \wedge b)$.

Уравнения (12.1) и (12.5) часто называют ► **аксиомами Колмогорова** в честь математика Андрея Колмогорова, показавшего, как построить остальную часть теории вероятностей на этом простом фундаменте и как справиться с трудностями,

вызванными непрерывными переменными.⁴ Хотя уравнение (12.2) имеет определенную особенность, уравнение (12.5) показывает, что аксиомы действительно ограничивают степень уверенности, которую агент может иметь в отношении логически связанных высказываний. Это аналогично тому факту, что логический агент не может одновременно быть уверен в высказываниях A , B и $\neg(A \wedge B)$, так как не существует возможного мира, в котором они все три одновременно являются истинными. Однако при использовании вероятностей высказывания относятся не к миру непосредственно, а к собственному состоянию знания агента. Почему же тогда агент не может придерживаться следующего множества убеждений (даже если они нарушают аксиомы Колмогорова)?

$$P(a) = 0,4 \quad P(b) = 0,3 \quad P(a \wedge b) = 0,0 \quad P(a \vee b) = 0,8 \quad (12.6)$$

Такого рода вопрос был предметом жаростных дебатов, продолжавшихся в течение десятилетий между теми, кто отстаивал допустимость использования вероятностей как единственной обоснованной формы оценки степеней уверенности, и теми, кто отстаивал альтернативные подходы.

Один из аргументов в пользу аксиом вероятностей, впервые сформулированный в 1931 году Бруно де Финетти ([554], 1983), заключается в следующем. Если агент имеет некоторую степень уверенности в истинности высказывания a , то он должен быть способен сформулировать оценку того, в какой степени он безразличен к ставке за или против высказывания a .⁵ Можно рассматривать подобную ситуацию как игру между двумя агентами: агент 1 утверждает: “Моя степень уверенности в истинности события a равна 0,4”. Затем агент 2 вправе выбрать, будет ли он делать ставку за или против высказывания a , выбирая ставки, совместимые с заявленной степенью уверенности. То есть агент 2 может решить сделать ставку на то, что событие a произойдет, поставив 6 долл. против 4 долл. агента 1. Или же агент 2 может сделать ставку на то, что будет иметь место событие $\neg a$, поставив 4 долл. против 6 агента 1. Когда исход события a станет известен, тот, кто оказался прав, заберет деньги. Если степень уверенности агента недостаточно точно отражает состояние мира, можно рассчитывать на то, что в долгосрочной перспективе он будет проигрывать деньги агенту-противнику, убеждения которого более точно отражают его состояние.

Теорема де Финетти относится не к выбору правильных значений для отдельных вероятностей, а к выбору значений вероятностей логически связанных

⁴ Эти трудности включают множество Витали, четко определенного измеримого подмножества в интервале $[0, 1]$ с неопределенной неизмеримой длиной.

⁵ Можно возразить, что предпочтения агента применительно к балансам разных ставок являются таковыми, что возможность потерять 1 долл. не уравнивается равной возможностью выиграть 1 долл. Один из возможных ответов на подобное возражение состоит в том, что суммы ставок должны быть достаточно малыми для того, чтобы можно было избежать данной проблемы. Анализ, проведенный Сэвджем ([1984], 1954), позволяет полностью исключить из рассмотрения эту проблему.

высказываний. ➔ Если агент 1 руководствуется множеством степеней уверенности, нарушающим аксиомы теории вероятностей, то всегда существует комбинация ставок агента 2, гарантирующая, что агент 1 будет терять деньги при каждой ставке. Например, предположим, что агент 1 руководствуется множеством степеней уверенности, приведенным в уравнении (12.6). На рис. 12.2 показано, что если агент 2 решит ставить 4 долл. на a , 3 долл. — на b и 2 долл. — на $\neg(a \vee b)$, то агент 1 всегда будет терять деньги, независимо от исходов для a и b . Из теоремы де Финетти следует, что ни один рациональный агент не может иметь убеждений, нарушающих аксиомы вероятности.

Высказывание	Степень уверенности агента 1	Ставка агента 2	Ставка агента 1	Результаты для агента 1			
				a, b	$a, \neg b$	$\neg a, b$	$\neg a, \neg b$
a	0,4	4 на a	6 на $\neg a$	-6	-6	4	4
b	0,3	3 на b	7 на $\neg b$	-7	3	-7	3
$a \vee b$	0,8	2 на $\neg(a \vee b)$	8 на $a \vee b$	2	2	2	-8
				-11	-1	-1	-1

Рис. 12.2. Поскольку агент 1 придерживается несогласованных убеждений, агент 2 может подобрать множество из трех ставок, гарантирующих постоянный проигрыш для агента 1, независимо от исходов для a и b

Одно общее возражение в отношении теоремы де Финетти состоит в том, что эта игра со ставками является довольно надуманной. Например, что будет, если один из игроков откажется делать ставку? Закончится ли на этом спор? Ответ на данный вопрос состоит в том, что эта игра со ставками представляет собой абстрактную модель для ситуации принятия решений, в которую любой агент неизбежно вовлечен в любой момент. Каждое действие (включая бездействие) — это своего рода ставка, а каждый исход может рассматриваться как положительное или отрицательное вознаграждение за эту ставку. Отказ делать ставку подобен отказу позволить времени двигаться.

В пользу применения вероятностей были выдвинуты и другие весомые философские аргументы, из которых наиболее заметными можно считать работы Кокса ([489], 1946), Карнапа ([374], 1950) и Джейнса (2003). В каждой из них предлагается множество аксиом для рассуждений со степенями доверия: отсутствие противоречий, соответствие положениям обычной логики (например, если степень доверия к A возрастает, то степень доверия к $\neg A$ должна уменьшаться) и т.д. Единственная спорная аксиома состоит в том, что степени доверия должны быть представлены числами или по крайней мере вести себя, как числа, например обладать свойством транзитивности (если степень доверия к A больше, чем степень доверия к B , которая больше, чем степень доверия к C , то степень доверия к A должна

быть больше, чем к C) и свойством сравнимости (степень доверия к A должна быть либо равна, либо больше, либо меньше, чем степень доверия к B). Можно доказать, что применение вероятностей является единственным подходом, удовлетворяющим все эти аксиомы.

Но мир таков, каков он есть, и практические свидетельства иногда оказываются более весомыми, чем доказательства. Успех систем формирования рассуждений, основанных на теории вероятностей, оказался гораздо более эффективным аргументом в пользу пересмотра многих взглядов, чем любая философская аргументация. В следующем разделе показано, как приведенные выше аксиомы можно применить к логическому выводу.

12.3. Логический вывод с использованием полных совместных распределений

В этом разделе описывается простой метод ► **вероятностного вывода**, т.е. вычисления апостериорных вероятностей для высказываний, сформулированных как ► **запросы** на основании наблюдаемых свидетельств. В качестве “базы знаний”, из которой можно будет вывести ответы на все запросы, мы будем использовать полное совместное распределение. По ходу дела также будет представлено несколько полезных методов манипулирования уравнениями, включающих вероятности.

Начнем с очень простого примера — проблемной области, состоящей только из трех булевых переменных, *Toothache*, *Cavity* и *Catch* (*щипцы*) (неприятные ощущения от захвата зуба стальными стоматологическими щипцами все еще свежи в памяти автора). Полное совместное распределение этих переменных представляет собой таблицу размером $2 \times 2 \times 2$, представленную на рис. 12.3.

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	0,108	0,012	0,072	0,008
\neg <i>cavity</i>	0,016	0,064	0,144	0,576

Рис. 12.3. Полное совместное распределение для мира *Toothache*, *Cavity*, *Catch*

Обратите внимание, что вероятности в этом совместном распределении в сумме составляют 1, что и требуется согласно аксиомам вероятности. Также обратите внимание, что уравнение (12.2) предоставляет прямой способ вычисления вероятности любого высказывания, простого или сложного: достаточно определить те возможные миры, в которых данное высказывание является истинным, а затем сложить их вероятности. Например, имеется шесть возможных миров, в которых высказывание $cavity \vee toothache$ является истинным:

$$P(\text{cavity} \vee \text{toothache}) = 0,108 + 0,012 + 0,072 + 0,008 + 0,016 + 0,064 = 0,28.$$

Одна из задач, которые встречаются особенно часто, состоит в том, чтобы извлечь из подобной таблицы распределение вероятностей по некоторому подмножеству переменных или по одной переменной. Например, складывая элементы первой строки на рис. 12.3, получим безусловную или **► маргинальную вероятность**⁶ события *cavity*:

$$P(\text{cavity}) = 0,108 + 0,012 + 0,072 + 0,008 = 0,2.$$

Этот процесс называется **► маргинализацией** или **исключением из суммы**, поскольку суммирование вероятностей для *каждого* возможного значения других переменных исключает их из уравнения. Можно записать следующее общее правило маргинализации для любых множеств переменных **Y** и **Z**:

$$P(\mathbf{Y}) = \sum_{\mathbf{z}} P(\mathbf{Y}, \mathbf{Z} = \mathbf{z}), \quad (12.7)$$

где $\sum_{\mathbf{z}}$ — сумма по всем возможным комбинациям значений множества переменных **Z**. Как обычно, в этом уравнении мы можем сократить $P(\mathbf{Y}, \mathbf{Z} = \mathbf{z})$ до $P(\mathbf{Y}, \mathbf{z})$. Например, для переменной *Cavity* уравнение (12.7) соответствует следующему уравнению:

$$\begin{aligned} P(\text{Cavity}) &= P(\text{Cavity}, \text{toothache}, \text{catch}) + P(\text{Cavity}, \text{toothache}, \neg\text{catch}) + \\ &+ P(\text{Cavity}, \neg\text{toothache}, \text{catch}) + P(\text{Cavity}, \neg\text{toothache}, \neg\text{catch}) = \\ &= \langle 0,108; 0,016 \rangle + \langle 0,012; 0,064 \rangle + \langle 0,072; 0,144 \rangle + \langle 0,008; 0,576 \rangle = \\ &= \langle 0,2; 0,8 \rangle. \end{aligned}$$

Используя правило умножения вероятностей (уравнение (12.4)), можно заменить $P(\mathbf{Y}, \mathbf{z})$ в уравнении (12.7) на $P(\mathbf{Y}|\mathbf{z})P(\mathbf{z})$, получив правило, называемое **► правилом обусловливания**:

$$P(\mathbf{Y}) = \sum_{\mathbf{z}} P(\mathbf{Y}|\mathbf{z})P(\mathbf{z}). \quad (12.8)$$

Как оказалось, правила маргинализации и обусловливания являются очень полезными правилами для всех видов логических выводов, включающих вероятностные выражения.

В большинстве случаев нас будет интересовать задача вычисления *условных* (апостериорных) вероятностей некоторых переменных при наличии свидетельств, касающихся других переменных. Условные вероятности можно найти, вначале воспользовавшись уравнением (12.3) для получения выражения в терминах безусловных вероятностей, а затем рассчитав это выражение на основании полного

⁶ Эта вероятность получила такое название, поскольку страховщики имеют общую привычку записывать суммы наблюдаемых частот событий на полях (*margin*) таблиц страхования.

совместного распределения. Например, ниже показано, как можно вычислить вероятность наличия полости в зубе, получив свидетельство о наличии зубной боли:

$$\begin{aligned} P(\text{cavity} | \text{toothache}) &= \frac{P(\text{cavity} \wedge \text{toothache})}{P(\text{toothache})} = \\ &= \frac{0,108 + 0,012}{0,108 + 0,012 + 0,016 + 0,064} = 0,6. \end{aligned}$$

Просто для проверки можно также рассчитать вероятность того, что у пациента нет полости в зубе, если у него наблюдается зубная боль:

$$\begin{aligned} P(\neg \text{cavity} | \text{toothache}) &= \frac{P(\neg \text{cavity} \wedge \text{toothache})}{P(\text{toothache})} = \\ &= \frac{0,016 + 0,064}{0,108 + 0,012 + 0,016 + 0,064} = 0,4. \end{aligned}$$

Эти два значения в сумме дают 1,0, как и должно быть. Обратите внимание на присутствие термина $P(\text{toothache})$ в знаменателе обоих этих вычислений. Если бы переменная *Cavity* имела более двух значений, этот терм присутствовал бы в знаменателе для всех них. Фактически для распределения $P(\text{Cavity} | \text{toothache})$ его можно рассматривать как константу **нормализации**, гарантирующую, что полученные вероятности в сумме составят 1. Во всех главах, в которых речь будет идти о вероятностях, для обозначения таких констант мы будем использовать символ α . Применив это обозначение, можно записать два предыдущих уравнения как одно:

$$\begin{aligned} \mathbf{P}(\text{Cavity} | \text{toothache}) &= \alpha \mathbf{P}(\text{Cavity}, \text{toothache}) = \\ &= \alpha [\mathbf{P}(\text{Cavity}, \text{toothache}, \text{catch}) + \mathbf{P}(\text{Cavity}, \text{toothache}, \neg \text{catch})] = \\ &= \alpha [(0,108; 0,016) + (0,012; 0,064)] = \alpha (0,12; 0,08) = (0,6; 0,4). \end{aligned}$$

Другими словами, мы можем вычислить $\mathbf{P}(\text{Cavity} | \text{toothache})$, даже не зная значения $P(\text{toothache})$! Забыв на время о множителе $1/P(\text{toothache})$, мы суммируем значения для *cavity* и $\neg \text{cavity}$, получив значения 0,12 и 0,08. Это правильная относительная пропорция, но в сумме они не дают 1, поэтому мы нормализуем эти значения делением каждого из них на $0,12 + 0,08$, получив в результате истинные вероятности 0,6 и 0,4. Нормализация оказывается полезным упрощением во многих вероятностных расчетах, позволяя как сделать вычисления проще, так и выполнить расчеты, когда некоторые оценки вероятности (например, $P(\text{toothache})$) недоступны.

На основании приведенного выше примера можно описать общую процедуру вероятностного вывода. Начнем со случая, когда запрос касается только одной переменной, X (например, *Cavity*). Пусть \mathbf{E} будет списком переменных свидетельства (в нашем примере это только *Toothache*), \mathbf{e} будет списком наблюдаемых значений

этих переменных, а \mathbf{Y} будет представлять оставшиеся ненаблюдаемые переменные (в нашем примере это только *Catch*). Запрос будет иметь вид $\mathbf{P}(X|\mathbf{e})$ и может быть вычислен следующим образом:

$$\mathbf{P}(X|\mathbf{e}) = \alpha \mathbf{P}(X, \mathbf{e}) = \alpha \sum_{\mathbf{y}} \mathbf{P}(X, \mathbf{e}, \mathbf{y}), \quad (12.9)$$

где суммирование осуществляется по всем возможным \mathbf{y} (т.е. по всем возможным комбинациям значений ненаблюдаемых переменных \mathbf{Y}). Обратите внимание, что взятые вместе переменные X , \mathbf{E} и \mathbf{Y} образуют полное множество переменных для данной проблемной области, поэтому $\mathbf{P}(X, \mathbf{e}, \mathbf{y})$ представляет собой просто подмножество вероятностей из полного совместного распределения.

При наличии полного совместного распределения, с которым можно работать, уравнение (12.9) позволяет получить ответы на вероятностные запросы в отношении дискретных переменных. Однако оно недостаточно хорошо масштабируется, поскольку в проблемной области с n булевыми переменными требуется входная таблица размером $O(2^n)$, обработка которой потребует времени $O(2^n)$. В реальных задачах вполне могут присутствовать сотни случайных переменных, и тогда оценка $O(2^n)$ для потребности в памяти и времени расчетов далеко выходит за пределы возможного: $2^{100} \approx 10^{30}$! И проблема здесь не только в объемах памяти и времени расчетов — еще более серьезная проблема состоит в том, что потребуется отдельно оценить на реальных примерах каждую из 10^{30} вероятностей, что делает объем необходимых экспериментальных данных просто астрономическим.

По этим причинам полное совместное распределение в табличной форме редко воспринимается как практический инструмент для построения систем рассуждений. На самом деле его, скорее, следует рассматривать как теоретическую основу, на которой можно строить более эффективные подходы, — как таблицы истинности являются теоретической основой для более практичных алгоритмов, таких как алгоритм DPLL, рассматривавшийся в главе 7. В оставшейся части этой главы будут представлены некоторые основные идеи, необходимые для подготовки к разработке реально осуществимых систем, описанных в главе 13.

12.4. Независимость

Давайте расширим полное совместное распределение, приведенное на рис. 12.2, добавив в него четвертую переменную, *Weather*. В результате полное совместное распределение будет иметь вид $\mathbf{P}(\textit{Toothache}, \textit{Catch}, \textit{Cavity}, \textit{Weather})$ и представлять собой таблицу из $2 \times 2 \times 2 \times 4 = 32$ элементов (переменная *Weather* имеет четыре значения). Это распределение содержит четыре “варианта” таблицы, представленной на рис. 12.3, по одному на каждый вид погоды. Возникает вопрос: какую связь эти варианты имеют друг с другом и с первоначальной таблицей, построенной при наличии трех переменных? Как связаны друг с другом значения

$P(\text{toothache}, \text{catch}, \text{cavity}, \text{cloud})$ и значения $P(\text{toothache}, \text{catch}, \text{cavity})$? Для получения ответа можно воспользоваться правилом умножения вероятностей (уравнение (12.4)):

$$P(\text{toothache}, \text{catch}, \text{cavity}, \text{cloud}) = P(\text{cloud} | \text{toothache}, \text{catch}, \text{cavity})P(\text{toothache}, \text{catch}, \text{cavity}).$$

Но человек, не верящий в возможность вмешательства свыше, едва ли сможет представить, что чьи-то проблемы с зубами способны повлиять на погоду. Поэтому следующее утверждение кажется вполне разумным:

$$P(\text{cloud} | \text{toothache}, \text{catch}, \text{cavity}) = P(\text{cloud}). \tag{12.10}$$

Из этого можно вывести следующее:

$$P(\text{toothache}, \text{catch}, \text{cavity}, \text{cloud}) = P(\text{cloud})P(\text{toothache}, \text{catch}, \text{cavity}).$$

Аналогичное уравнение существует для *каждого* элемента в распределении $\mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather})$. В действительности можно даже записать такое общее уравнение:

$$\mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather}) = \mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity})\mathbf{P}(\text{Weather}).$$

Следовательно, 32-элементная таблица для четырех переменных может быть образована из одной 8-элементной таблицы и одной 4-элементной. Подобная декомпозиция схематически показана на рис. 12.4, а.

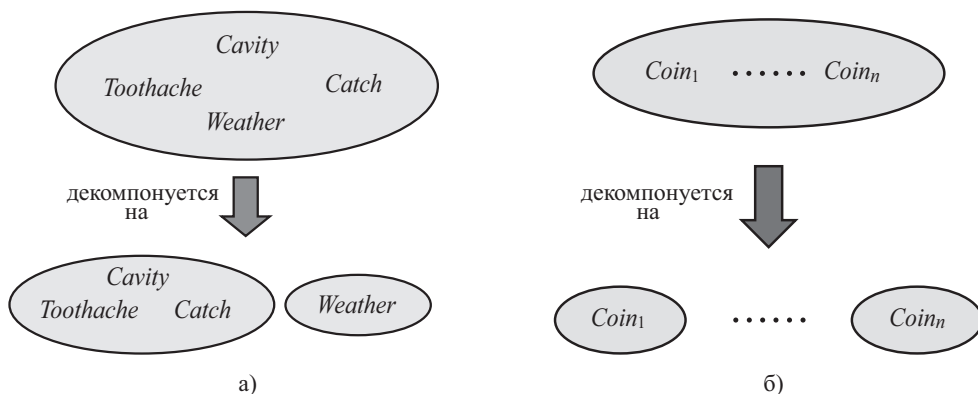


Рис. 12.4. Два примера разбиения большого совместного распределения на меньшие распределения с использованием свойства абсолютной независимости. **а)** Погода и проблемы с зубами независимы друг от друга. **б)** Броски монеты независимы друг от друга

Свойство, которое использовалось в уравнении (12.10), называется ► **независимостью** (а также **маргинальной независимостью** или **абсолютной независимостью**). В частности, погода независима от чьих-то проблем с зубами. Независимость между высказываниями a и b может быть представлена следующим образом:

$$P(a|b) = P(a) \quad \text{или} \quad P(b|a) = P(b) \quad \text{или} \quad P(a \wedge b) = P(a)P(b). \quad (12.11)$$

Все эти варианты записи эквивалентны (упражнение 12.15). Свойство независимости между переменными X и Y можно сформулировать следующим образом (все эти варианты записи также эквивалентны):

$$\mathbf{P}(X|Y) = \mathbf{P}(X) \quad \text{или} \quad \mathbf{P}(Y|X) = \mathbf{P}(Y) \quad \text{или} \quad \mathbf{P}(X, Y) = \mathbf{P}(X)\mathbf{P}(Y).$$

Утверждения о независимости обычно основаны на знаниях о проблемной области. Как показывает пример с зубной болью и погодой, они позволяют существенно сократить объем информации, необходимой для построения полного совместного распределения. Если все множество переменных может быть разделено на независимые подмножества, то полное совместное распределение может быть *факторизовано* на отдельные совместные распределения, заданные на этих подмножествах. Например, совместное распределение результатов n независимых бросков монеты $\mathbf{P}(C_1, \dots, C_n)$ включает 2^n элементов, но может быть представлено как произведение n распределений $\mathbf{P}(C_i)$ с одной переменной. С практической точки зрения независимость стоматологии и метеорологии является благоприятным фактором, поскольку в противном случае стоматологам могли бы потребоваться глубокие знания в области метеорологии, и наоборот.

Поэтому любые утверждения о независимости, если они имеются, позволяют сократить размеры представления проблемной области и уменьшить сложность проблемы логического вывода. К сожалению, чистое разделение полных множеств переменных по признаку независимости встречается редко. Если между двумя переменными существует хоть какая-то связь, пусть даже косвенная, свойство независимости уже не соблюдается. Более того, даже независимые подмножества могут оказаться достаточно большими; например, в области стоматологии могут рассматриваться десятки заболеваний и сотни симптомов, причем все они взаимосвязаны друг с другом. Чтобы справиться с такими задачами, потребуются методы, более тонкие, чем прямолинейная концепция независимости.

12.5. Правило Байеса и его использование

В разделе 12.2.1 было определено правило умножения вероятностей (уравнение (12.4)). На самом деле это правило можно записать в двух формах:

$$P(a \wedge b) = P(a|b)P(b) \quad \text{и} \quad P(a \wedge b) = P(b|a)P(a).$$

Приравняв правые части этих двух уравнений и разделив их на $P(a)$, получим

$$P(a \wedge b) = \frac{P(a|b)P(b)}{P(a)}. \quad (12.12)$$

Это уравнение известно как ► **правило Байеса** (закон Байеса, теорема Байеса). Это простое уравнение лежит в основе всех современных систем искусственного интеллекта для вероятностного вывода.

Более общий случай правила Байеса для многозначных переменных можно записать в нотации **P** следующим образом:

$$\mathbf{P}(Y|X) = \frac{\mathbf{P}(X|Y)\mathbf{P}(Y)}{\mathbf{P}(X)}.$$

Как и прежде, это уравнение также следует рассматривать как представляющее множество уравнений, в каждом из которых рассматриваются конкретные значения переменных. Время от времени нам также придется использовать более общую версию, которая обусловлена некоторым фоновым свидетельством **e**:

$$\mathbf{P}(Y|X, \mathbf{e}) = \frac{\mathbf{P}(X|Y, \mathbf{e})\mathbf{P}(Y|\mathbf{e})}{\mathbf{P}(X|\mathbf{e})}. \quad (12.13)$$

12.5.1. Применение правила Байеса: простой случай

На первый взгляд, правило Байеса не кажется очень полезным. Оно позволяет вычислить единственный терм $P(b|a)$ на основании трех термов, $P(a|b)$, $P(b)$ и $P(a)$. Хотя это и похоже на шаг вперед и два шага назад, тем не менее правило Байеса находит очень широкое практическое применение, поскольку во многих случаях имеются хорошие оценки вероятностей для этих трех элементов и нужно вычислить четвертый. Часто воспринимаемые данные указывают на *результат* (*effect*), вызываемый какой-то неизвестной *причиной* (*cause*), и нам необходимо определить эту причину. В таких случаях правило Байеса принимает следующий вид:

$$P(\text{cause} | \text{effect}) = \frac{P(\text{effect} | \text{cause})P(\text{cause})}{P(\text{effect})}.$$

Условная вероятность $P(\text{effect} | \text{cause})$ предоставляет количественную оценку взаимосвязи в ► **причинном** направлении, тогда как вероятность $P(\text{cause} | \text{effect})$ описывает ► **диагностическое** направление. В такой задаче, как определение медицинского диагноза, мы часто имеем условные вероятности причинно-следственных связей. Врач знает диагностические вероятности $P(\text{symptoms} | \text{disease})$ и хочет поставить диагноз $P(\text{disease} | \text{symptoms})$.

Например, врач знает, что менингит часто вызывает у пациента снижение подвижности шеи, — предположим, это наблюдается в 70% случаев. Врач также знает некоторые безусловные факты: априорная вероятность того, что у пациента менингит, равна 1/50 000, а априорная вероятность того, что у пациента будет сни-

жена подвижность шеи, равна 1%. Пусть s — высказывание, утверждающее, что у пациента снижена подвижность шеи, а m — высказывание, утверждающее, что у пациента менингит. Тогда можно записать следующее:

$$\begin{aligned} P(s | m) &= 0,7 \\ P(m) &= 1/50\,000 \\ P(s) &= 0,01 \\ P(m | s) &= \frac{P(s | m)P(m)}{P(s)} = \frac{0,7 \times 1/50\,000}{0,01} = 0,0014. \end{aligned} \quad (12.14)$$

Таким образом, можно ожидать, что только у 0,14% пациентов со сниженной подвижностью шеи будет менингит. Обратите внимание, что даже если снижение подвижности шеи является весьма надежным свидетельством наличия менингита (с вероятностью 0,7), сама вероятность наличия менингита у пациента остается очень низкой. Это связано с тем, что априорная вероятность симптома снижения подвижности шеи (по любой причине) намного выше в сравнении с вероятностью менингита.

В разделе 12.3 был описан процесс, посредством которого можно избежать необходимости оценки априорной вероятности свидетельства (в данном случае — $P(s)$), вычислив вместо этого апостериорную вероятность для каждого значения переменной запроса (в данном случае — m и $\neg m$), а затем нормализовав результаты. Аналогичный процесс можно применять и при использовании правила Байеса. Итак, мы имеем

$$\mathbf{P}(M | s) = \alpha \langle P(s | m)P(m), P(s | \neg m)P(\neg m) \rangle.$$

Следовательно, чтобы воспользоваться этим подходом, необходимо вместо $P(s)$ вычислить значение $P(s | \neg m)$. К сожалению, бесплатных пирожных не бывает, — иногда это упрощает задачу, а иногда усложняет. Общая форма правила Байеса с нормализацией будет следующей:

$$\mathbf{P}(Y | X) = \alpha \mathbf{P}(X | Y)\mathbf{P}(Y), \quad (12.15)$$

где α — константа нормализации, необходимая для того, чтобы сумма элементов в распределении $\mathbf{P}(Y | X)$ была равна 1.

Один из очевидных вопросов, касающихся правила Байеса, состоит в том, почему доступной может оказаться условная вероятность, реализуемая только в одном направлении, но не в другом. В проблемной области лечения менингита врач, возможно, знает, что при наличии симптома ограничения подвижности шеи менингит будет причиной в 1 из 5000 случаев. А это означает, что у врача имеется количественная информация в **диагностическом** направлении, от симптома к причине. Такому врачу использовать правило Байеса не требуется.

К сожалению, ➔ *знания в диагностическом направлении на практике встречаются намного реже, чем знания в причинном направлении*. В случае внезапной эпидемии менингита априорная вероятность этого заболевания, $P(m)$, повышается. Врач,

который вывел диагностическую вероятность $P(m | s)$ непосредственно из статистических наблюдений за пациентами перед эпидемией, не будет иметь представления о том, как обновить это значение после ее начала, тогда как врач, вычисляющий значение $P(m | s)$ из других трех значений, быстро обнаружит, что значение $P(m | s)$ должно увеличиваться пропорционально $P(m)$. Еще более важно то, что причинная информация $P(m | s)$ остается незатронутой эпидемией, поскольку она просто отражает, в чем выражается воздействие менингита на пациента. Использование прямых причинных знаний такого рода или знаний, основанных на модели, позволяет достичь надежности, которая крайне важна при создании вероятностных систем, применимых в реальном мире.

12.5.2. Использование правила Байеса: комбинирование свидетельств

Выше было показано, что правило Байеса может применяться для получения ответов на вероятностные запросы, в которых учтено условие, составляющее одно из свидетельств, например ограниченная подвижность шеи. В частности, было показано, что вероятностная информация часто доступна в форме $P(\text{effect} | \text{cause})$, где *effect* — результат, а *cause* — причина. А что произойдет, если свидетельств два или больше? Например, какой вывод может сделать стоматолог, если его пугающие стальные щипцы сомкнулись на больном зубе пациента? Если известно полное совместное распределение (см. рис. 12.2), можно сразу же отыскать ответ:

$$\mathbf{P}(\text{Cavity} | \text{toothache} \wedge \text{catch}) = \alpha \langle 0,108; 0,016 \rangle \approx \langle 0,871; 0,129 \rangle.$$

Но нам уже известно, что такой подход не масштабируется на большее количество переменных. Можно попробовать воспользоваться правилом Байеса для переформулировки этой задачи:

$$\begin{aligned} \mathbf{P}(\text{Cavity} | \text{toothache} \wedge \text{catch}) &= \\ &= \alpha \mathbf{P}(\text{toothache} \wedge \text{catch} | \text{Cavity}) \mathbf{P}(\text{Cavity}). \end{aligned} \tag{12.16}$$

Чтобы получить ответ на запрос в такой формулировке, необходимо знать условные вероятности конъюнкции $\text{toothache} \wedge \text{catch}$ для каждого значения *Cavity*. Это может быть легко осуществимо, если речь идет только о двух переменных свидетельства, но такой подход вновь становится источником затруднений при его масштабировании. Если имеется n возможных переменных свидетельства (рентгеновский снимок, гигиена полости рта и т.д.), то количество возможных комбинаций наблюдаемых значений, для которых необходимо будет знать условные вероятности, составит $O(2^n)$. Это не лучше, чем использование полного совместного распределения.

Чтобы улучшить ситуацию, необходимо найти некоторые дополнительные утверждения о рассматриваемой проблемной области, позволяющие упростить применяемые выражения. Понятие **независимости**, введенное в разделе 12.4, дает

ключ к этому решению, но требует уточнения. Было бы прекрасно, если бы переменные *Toothache* и *Catch* были независимыми, но они таковыми не являются: если зубной врач захватывает зуб щипцами, то он делает это, вероятно, потому, что в этом зубе есть полость, и именно наличие этой полости вызывает боль. Однако эти переменные *действительно* являются независимыми, *когда речь идет о наличии или отсутствии полости*. Причиной в каждом случае действительно является наличие полости в зубе, но ни одна из этих переменных не оказывает непосредственного влияния на другую: зубную боль определяет состояние нервов в зубе, тогда как точность наложения инструмента зависит прежде всего от навыков стоматолога, к которым зубная боль не имеет никакого отношения.⁷ Математически это свойство записывается следующим образом:

$$\mathbf{P}(\textit{toothache} \wedge \textit{catch} \mid \textit{Cavity}) = \mathbf{P}(\textit{toothache} \mid \textit{Cavity})\mathbf{P}(\textit{catch} \mid \textit{Cavity}). \quad (12.17)$$

В этом уравнении выражена ► **условная независимость** переменных *toothache* и *catch*, если дана вероятность *Cavity*. Соответствующее выражение можно вставить в уравнение (12.16) с целью определения вероятности наличия полости:

$$\begin{aligned} \mathbf{P}(\textit{Cavity} \mid \textit{toothache} \wedge \textit{catch}) &= \\ &= \alpha \mathbf{P}(\textit{toothache} \mid \textit{Cavity}) \mathbf{P}(\textit{catch} \mid \textit{Cavity}) \mathbf{P}(\textit{Cavity}). \end{aligned} \quad (12.18)$$

Теперь требования к наличию информации становятся такими же, как и при вероятностном выводе с использованием каждого свидетельства в отдельности: необходимо знать априорную вероятность $\mathbf{P}(\textit{Cavity})$ для переменной запроса и условную (апостериорную) вероятность каждого результата, если дана его причина.

Общее определение **условной независимости** двух переменных, *X* и *Y*, если дана третья переменная, *Z*, выражается следующей формулой:

$$\mathbf{P}(X, Y \mid Z) = \mathbf{P}(X \mid Z) \mathbf{P}(Y \mid Z).$$

Например, в проблемной области стоматологии кажется вполне резонным применить утверждение об условной независимости переменных *Toothache* и *Catch*, если дана вероятность *Cavity*:

$$\begin{aligned} \mathbf{P}(\textit{Toothache}, \textit{Catch} \mid \textit{Cavity}) &= \\ &= \mathbf{P}(\textit{Toothache} \mid \textit{Cavity}) \mathbf{P}(\textit{Catch} \mid \textit{Cavity}). \end{aligned} \quad (12.19)$$

Обратите внимание, что это утверждение несколько строже по сравнению с уравнением (12.17), в котором утверждается независимость только для конкретных значений *Toothache* и *Catch*. Если же воспользоваться свойством абсолютной независимости (уравнение (12.11)), то получим следующие эквивалентные формы, которыми также можно пользоваться (упражнение 12.21):

$$\mathbf{P}(X \mid Y, Z) = \mathbf{P}(X \mid Z) \quad \text{и} \quad \mathbf{P}(Y \mid X, Z) = \mathbf{P}(Y \mid Z).$$

⁷ Предполагается, что пациент и стоматолог — разные люди.

В разделе 12.4 было показано, что утверждения при наличии абсолютной независимости позволяют выполнять декомпозицию полного совместного распределения на гораздо более мелкие распределения. Как оказалось, аналогичную декомпозицию допускают и утверждения при наличии условной независимости. Например, для утверждения в уравнении (12.19) декомпозицию можно вывести следующим образом:

$$\begin{aligned} P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}) &= \\ &= P(\textit{Toothache}, \textit{Catch} \mid \textit{Cavity}) P(\textit{Cavity}) = && \text{(по правилу умножения)} \\ &= P(\textit{Toothache} \mid \textit{Cavity}) P(\textit{Catch} \mid \textit{Cavity}) P(\textit{Cavity}) && \text{(по уравнению (12.19)).} \end{aligned}$$

(То, что это уравнение действительно выполняется, читатель легко может проверить, обратившись к рис. 12.3.) Таким образом, большая исходная таблица теперь декомпонована на три меньшие таблицы. В исходной таблице было семь независимых значений. (Эта таблица включает $2^3 = 8$ значений, но их сумма должна быть равна 1, поэтому только 7 из них являются независимыми). Меньшие таблицы содержат в общей сложности $2 + 2 + 1 = 5$ независимых значений. (Распределение условных вероятностей, такое как $P(\textit{Toothache} \mid \textit{Cavity})$, включает две строки из двух значений, и в каждой строке их сумма равна 1, поэтому в данном случае есть только два независимых значения, а для априорного распределения, такого как $P(\textit{Cavity})$, существует только одно независимое значение.) Переход от 7 к 5 независимым значениям может показаться не таким уж большим достижением, но выигрыш может оказаться намного больше при увеличении количества симптомов.

В общем случае для n симптомов, все из которых являются условно независимыми при заданной вероятности \textit{Cavity} , размер представления растет как $O(n)$, а не $O(2^n)$. Это означает, что **→ утверждения об условной независимости могут обеспечить масштабирование вероятностных систем; более того, такие утверждения могут быть подкреплены данными намного проще по сравнению с утверждениями об абсолютной независимости.** С концептуальной точки зрения переменная \textit{Cavity} **▶ разделяет** переменные $\textit{Toothache}$ и \textit{Catch} , поскольку наличие полости в зубе является прямой причиной и зубной боли, и наложения щипцов на больной зуб. Разработка методов декомпозиции крупных вероятностных областей определения на слабо связанные подмножества с помощью свойства условной независимости стало одним из наиболее важных достижений в новейшей истории искусственного интеллекта.

12.6. Наивные байесовские модели

Приведенный выше пример из области стоматологии иллюстрирует часто встречающуюся ситуацию, в которой одна причина непосредственно влияет на целый ряд результатов, причем все эти результаты являются условно независимыми, когда дана эта причина. Полное совместное распределение может быть записано следующим образом:

$$P(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = P(\text{Cause}) \prod_i P(\text{Effect}_i | \text{Cause}). \quad (12.20)$$

Такое распределение вероятностей называется ► **наивной байесовской** моделью, “наивной” — потому что часто используется (как упрощающее допущение) в тех случаях, когда переменные “результата” *не являются* строго независимыми при данной переменной причины. (Наивную байесовскую модель иногда называют ► **байесовским классификатором**, и это не совсем корректное употребление термина побудило настоящих специалистов в области байесовских моделей называть ее не наивной, а **идиотской байесовской** моделью.) На практике наивные байесовские системы часто могут работать очень успешно, даже если предположение о независимости не является строго истинным.

Для использования наивной байесовской модели можно применить уравнение (12.20), чтобы получить вероятность причины с учетом некоторых наблюдаемых результатов. Обозначим наблюдаемые результаты как $\mathbf{E} = \mathbf{e}$, тогда как остальные переменные результата \mathbf{Y} являются ненаблюдаемыми. Далее можно применить стандартный метод логического вывода из совместного распределения (уравнение (12.9)):

$$P(\text{Cause} | \mathbf{e}) = \alpha \sum_{\mathbf{y}} P(\text{Cause}, \mathbf{e}, \mathbf{y}).$$

Из уравнения (12.20) получаем

$$\begin{aligned} P(\text{Cause} | \mathbf{e}) &= \alpha \sum_{\mathbf{y}} P(\text{Cause}) P(\mathbf{y} | \text{Cause}) \left(\prod_j P(e_j | \text{Cause}) \right) = \\ &= \alpha P(\text{Cause}) \left(\prod_j P(e_j | \text{Cause}) \right) \sum_{\mathbf{y}} P(\mathbf{y} | \text{Cause}) = \\ &= \alpha P(\text{Cause}) \prod_j P(e_j | \text{Cause}) \end{aligned} \quad (12.21)$$

Здесь последняя строка может быть выведена потому, что сумма по \mathbf{y} равна 1. Словами это уравнение можно интерпретировать так: для каждой возможной причины умножьте априорную вероятность причины на произведение условных вероятностей наблюдаемых результатов на данную причину, а затем нормализуйте результат. Время выполнения этих расчетов изменяется линейно по отношению к количеству наблюдаемых результатов и не зависит от количества ненаблюдаемых результатов (которое может быть очень большим в таких предметных областях, как медицина). В следующей главе будет показано, что это обычное явление в вероятностном выводе: переменные свидетельства, значения которых ненаблюдаемы, обычно “исчезают” из вычислений в полном составе.

12.6.1. Классификация текста с помощью наивной байесовской модели

Давайте посмотрим, как наивную байесовскую модель можно применить в задаче ► **классификации текстов**: дан некоторый текст, и необходимо установить, к какому из заранее определенного набора классов или категорий он относится. Здесь в качестве “причины” выступает переменная $Category$, а наличие или отсутствие в тексте определенных ключевых слов представлено переменными “результата” $HasWord_i$. Рассмотрим следующие два примера предложений, взятых из газетных статей.

1. В понедельник стоимость акций возросла — основные индексы прибавили 1%, поскольку сохраняется оптимизм в отношении сезона отчетности за первый квартал.
2. В понедельник проливные дожди по-прежнему охватывают большую часть восточного побережья, из-за чего в городе Нью-Йорк и других местах были сделаны предупреждения о возможности наводнения.

Наша задача заключается в том, чтобы отнести каждое предложение к некоторой $Category$ — основному разделу газет: новости (*news*), спорт (*sports*), бизнес (*business*), погода (*weather*) или развлечения (*entertainment*). Наивная байесовская модель включает априорные вероятности $P(Category)$ и условные вероятности $P(HasWord_i | Category)$. Для каждой категории c вероятность $P(Category = c)$ оценивается как доля документов, относящихся к этой категории, из числа всех ранее просмотренных документов. Например, если в 9% статей идет речь о погоде, то $P(Category = weather) = 0,09$. Аналогичным образом, вероятности $P(HasWord_i | Category)$ оцениваются как доля документов каждой категории, в которых присутствует слово i . Так, если примерно 37% статей о бизнесе содержат слово № 6, “акции”, то вероятности $P(HasWord_6 = true | Category = business)$ можно присвоить значение 0,37.⁸

Чтобы классифицировать новый документ, необходимо проверить, какие ключевые слова в нем присутствуют, а затем применить уравнение (12.21), чтобы получить распределение апостериорных вероятностей по категориям. Если необходимо указать только одну категорию, выбирается та, у которой будет наибольшая апостериорная вероятность. Обратите внимание, что в этой задаче каждая

⁸ Нужно проявлять осторожность, чтобы не присвоить нулевую вероятность словам, которые ранее не встречались в данной категории документов, поскольку нулевое значение уничтожит все остальные свидетельства в уравнении (12.21). То, что слово пока не встречалось, еще не означает, что оно никогда не встретится. Вместо этого следует резервировать небольшую часть распределения вероятностей для представления слов, “ранее не наблюдавшихся”. Читайте главу 20 для получения более подробной информации по этому вопросу в целом и раздел 23.1.4, в котором представлены конкретные примеры словесных моделей.

переменная результата является наблюдаемой, поскольку всегда можно с уверенностью сказать, присутствует данное слово в документе или нет.

В наивной байесовской модели предполагается, что слова в документах появляются независимо друг от друга с частотой, определяемой категорией документа. Это предположение о независимости, очевидно, нарушается на практике. Например, фраза “первый квартал” встречается в статьях о бизнесе (или спорте) чаще, чем можно было бы предположить путем умножения вероятностей для отдельных слов “первый” и “квартал”. Нарушение независимости обычно означает, что конечные апостериорные вероятности будут намного ближе к 1 или 0, чем они должны быть, — другими словами, такая модель будет проявлять излишнюю самоуверенность в своих предсказаниях. С другой стороны, даже с такими ошибками *рейтинг* возможных категорий часто оказывается весьма точным.

Наивные байесовские модели широко используются для определения языка и поиска документов, фильтрации спама и других задач классификации. Для решения таких задач, как медицинская диагностика, в которой фактические значения апостериорных вероятностей действительно имеют значение — например, при принятии решения, следует ли удалить аппендикс, — предпочтение отдают более сложным моделям, описываемым в следующей главе.

12.7. Очередное возвращение в мир вампуса

Комбинацию идей, изложенных в этой главе, можно применить для решения задачи выполнения вероятностных рассуждений в мире вампуса (полное описание мира вампуса приведено в главе 7). Неопределенность в мире вампуса возникает из-за того, что датчики агента предоставляют ему только частичную информацию о состоянии этого мира. Например, на рис. 12.5 показана ситуация, в которой каждый из трех не посещенных, но достижимых квадратов, [1,3], [2,2] и [3,1], может содержать яму. Чисто логический вывод не позволяет прийти к каким-либо заключениям о том, какой квадрат с наибольшей вероятностью окажется безопасным, поэтому логический агент может быть вынужден сделать случайный выбор. В этом разделе будет показано, что в подобной ситуации вероятностный агент может действовать гораздо успешнее, чем логический агент.

Наша цель — вычислить вероятность наличия ямы в каждом из этих трех квадратов. (В данном конкретном примере присутствие в них вампуса и золота игнорируется.) Относящиеся к этой задаче свойства мира вампуса включают, во-первых, то, что наличие ямы вызывает ощущение ветерка во всех соседних квадратах, и во-вторых то, что в каждом квадрате, отличном от [1,1], вероятность наличия в нем ямы равна 0,2. На первом этапе определяем множество необходимых случайных переменных, как показано ниже.

- Как и в случае логики высказываний, нам потребуется по одной булевой переменной P_{ij} для каждого квадрата; она будет принимать значение *true* тогда и только тогда, когда квадрат $[i,j]$ действительно содержит яму.
- Также нам потребуются булевы переменные B_{ij} , принимающие значение *true* тогда и только тогда, когда в квадрате $[i,j]$ ощущается ветерок. Из общего числа этих переменных нам достаточно будет рассмотреть только те, которые относятся к наблюдаемым квадратам, в данном случае — $[1,1]$, $[1,2]$ и $[2,1]$.

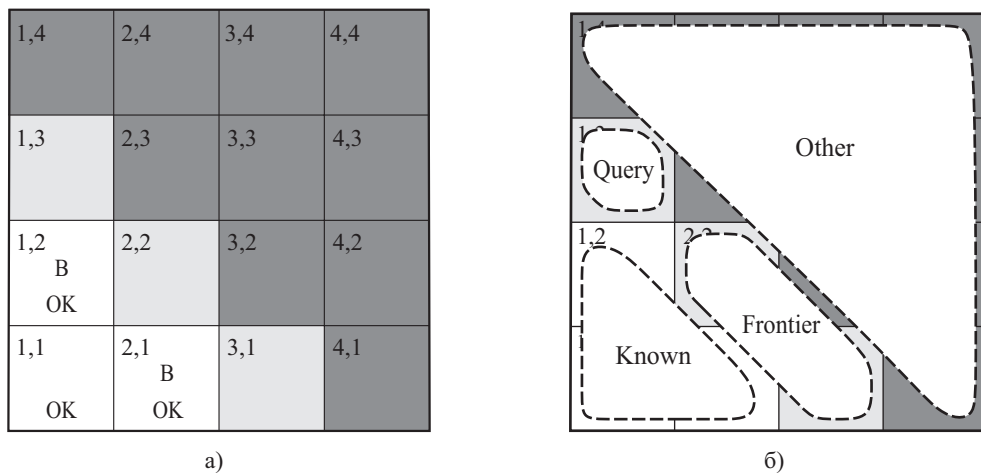


Рис. 12.5. а) После обнаружения ветерка как в квадрате $[1,2]$, так и в квадрате $[2,1]$ агент заходит в тупик — нет такого квадрата, который он мог бы обследовать без опасений. б) Распределение квадратов по категориям *Known* (известные), *Frontier* (периферийные) и *Other* (прочие) для формирования запроса (*Query*) в отношении квадрата $[1,3]$

На следующем этапе определяем полное совместное распределение $\mathbf{P}(P_{1,1}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1})$. Применив правило умножения вероятностей, получаем следующее:

$$\begin{aligned} \mathbf{P}(P_{1,1}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1}) &= \\ &= \mathbf{P}(B_{1,1}, B_{1,2}, B_{2,1} | P_{1,1}, \dots, P_{4,4}) \mathbf{P}(P_{1,1}, \dots, P_{4,4}). \end{aligned}$$

Эта декомпозиция позволяет очень легко определить, какими должны быть значения совместной вероятности. Первый терм представляет собой условную вероятность некоторой конфигурации данных о наличии ветерка, если дана конфигурация расположения ям; он принимает значение 1, если в квадратах, соседних с ямами, чувствуется ветерок, в противном случае принимает значение 0. Вторым термом является априорная вероятность конфигурации расположения ям. Каждый

квадрат может содержать яму с вероятностью 0,2, независимо от других квадратов, поэтому имеет место следующее:

$$\mathbf{P}(P_{1,1}, \dots, P_{4,4}) = \prod_{i,j=1,1}^{4,4} \mathbf{P}(P_{ij}). \quad (12.22)$$

Для конкретной конфигурации с точно n ямами эта вероятность равна $0,2^n \times 0,8^{16-n}$.

В ситуации, показанной на рис. 12.2, *а*, свидетельство состоит из наблюдаемого ветерка (или его отсутствия) в каждом посещенном квадрате в сочетании с тем фактом, что каждый такой квадрат не содержит ямы. Эти факты можно сокращенно представить как $b = \neg b_{1,1} \wedge b_{1,2} \wedge b_{2,1}$ и $known = \neg p_{1,1} \wedge \neg p_{1,2} \wedge \neg p_{2,1}$. Нас интересуют ответы на такие запросы, как $\mathbf{P}(P_{1,3} | known, b)$: насколько велика вероятность того, что квадрат [1,3] содержит яму, учитывая результаты всех наблюдений, сделанных до сих пор?

Чтобы получить ответ на этот запрос, можно использовать стандартный подход, основанный на уравнении (12.9), а именно — просто просуммировать элементы таблицы полного совместного распределения. Пусть *Unknown* (неизвестное) — это множество переменных P_{ij} для квадратов, отличных от уже проверенных квадратов и квадрата запроса [1,3]. В таком случае, следуя уравнению (12.9), получим:

$$\mathbf{P}(P_{1,3} | known, b) = \alpha \sum_{unknown} \mathbf{P}(P_{1,3}, known, b, unknown). \quad (12.23)$$

Полное совместное распределение вероятностей уже было определено, поэтому можно считать, что задача решена; точнее, осталось только выполнить вычисления. Количество неизвестных квадратов равно 12, следовательно, требуемая сумма состоит из $2^{12} = 4096$ термов. В общем случае количество термов в этой сумме растет экспоненциально в зависимости от количества квадратов.

Безусловно, напрашивается вопрос, а не являются ли другие квадраты не относящимися к делу? Как содержимое квадрата [4,4] может повлиять на наличие ямы в квадрате [1,3]? И действительно, эта догадка является приблизительно правильной, но ее необходимо уточнить. На самом деле здесь мы имеем в виду, что если бы мы знали значения переменных P для всех смежных квадратов, которые нас интересуют, то наличие (или отсутствие) ямы в других, более отдаленных, квадратах, уже не могло бы оказать влияния на нашу уверенность.

Пусть *Frontier* будет множеством переменных P (отличных от переменной запроса) всех тех квадратов, которые являются смежными с посещенными квадратами, — в нашем случае это квадраты [2,2] и [3,1]. Кроме того, пусть *Other* будет множеством переменных P для всех остальных неизвестных квадратов, — в нашем случае их имеется 10, как показано на рис. 12.5, *б*. С учетом этих определений $Unknown = Frontier \cup Other$. Ключевая догадка, приведенная выше, теперь может быть сформулирована следующим образом: наблюдения ветерка *условно*

независимы от других переменных, если даны известные переменные, переменные множества *Frontier* и переменная запроса. Чтобы воспользоваться этой идеей, необходимо преобразовать формулу запроса в такую форму, в которой данные о наличии ветерка становятся условно зависимыми от всех других переменных, а затем упростить полученное выражение с использованием утверждения об условной независимости:

$$\begin{aligned}
 \mathbf{P}(P_{1,3} | \textit{known}, b) &= \\
 &= \alpha \sum_{\textit{unknown}} \mathbf{P}(P_{1,3}, \textit{known}, b, \textit{unknown}) = && \text{(из уравнения (12.23))} \\
 &= \alpha \sum_{\textit{unknown}} \mathbf{P}(b | P_{1,3}, \textit{known}, \textit{unknown}) \mathbf{P}(P_{1,3}, \textit{known}, \textit{unknown}) = \\
 & && \text{(правило умножения вероятностей)} \\
 &= \alpha \sum_{\textit{frontier}} \sum_{\textit{other}} \mathbf{P}(b | \textit{known}, P_{1,3}, \textit{frontier}, \textit{other}) \mathbf{P}(P_{1,3}, \textit{known}, \textit{frontier}, \textit{other}) = \\
 &= \alpha \sum_{\textit{frontier}} \sum_{\textit{other}} \mathbf{P}(b | \textit{known}, P_{1,3}, \textit{frontier}) \mathbf{P}(P_{1,3}, \textit{known}, \textit{frontier}, \textit{other}),
 \end{aligned}$$

где на конечном этапе используется утверждение об условной независимости: переменная *b* не зависит от *других переменных*, если даны известные переменные, переменные множества *Frontier* и переменная запроса $P_{1,3}$. Теперь первый терм в выражении не зависит от переменных множества *Other*, поэтому операцию суммирования можно переместить внутрь выражения:

$$\begin{aligned}
 \mathbf{P}(P_{1,3} | \textit{known}, b) &= \\
 &= \alpha \sum_{\textit{frontier}} \mathbf{P}(b | \textit{known}, P_{1,3}, \textit{frontier}) \sum_{\textit{other}} \mathbf{P}(P_{1,3}, \textit{known}, \textit{frontier}, \textit{other}).
 \end{aligned}$$

Согласно утверждению о независимости, соответствующему приведенному в уравнении (12.22), терм априорной вероятности может быть факторизован, после чего все эти термы могут быть переупорядочены следующим образом:

$$\begin{aligned}
 \mathbf{P}(P_{1,3} | \textit{known}, b) &= \\
 &= \alpha \sum_{\textit{frontier}} \mathbf{P}(b | \textit{known}, P_{1,3}, \textit{frontier}) \sum_{\textit{other}} \mathbf{P}(P_{1,3}) P(\textit{known}) P(\textit{frontier}) P(\textit{other}) = \\
 &= \alpha P(\textit{known}) \mathbf{P}(P_{1,3}) \sum_{\textit{frontier}} \mathbf{P}(b | \textit{known}, P_{1,3}, \textit{frontier}) P(\textit{frontier}) \sum_{\textit{other}} P(\textit{other}) = \\
 &= \alpha' \mathbf{P}(P_{1,3}) \sum_{\textit{frontier}} \mathbf{P}(b | \textit{known}, P_{1,3}, \textit{frontier}) P(\textit{frontier}),
 \end{aligned}$$

где на последнем этапе постоянный терм $P(\textit{known})$ вводится в нормализующую константу на основании того факта, что выражение $\sum_{\textit{other}} P(\textit{other})$ равно 1.

Теперь в сумме по переменным множества *Frontier* $P_{2,2}$ и $P_{3,1}$ осталось только четыре термина. Использование свойств независимости и условной независимости позволило полностью исключить из рассмотрения все остальные квадраты.

Обратите внимание на то, что сумма вероятностей в выражении $\mathbf{P}(b | \textit{known}, P_{1,3}, \textit{frontier})$ равна 1, если данные наблюдений о наличии ветерка совместимы с другими переменными, — в противном случае она равна 0. Следовательно, для

каждого значения $P_{1,3}$ выполняется суммирование по *логическим моделям* для переменных в множестве *Frontier*, согласующимся с известными фактами (это можно сравнить с тем, как осуществлялся перебор моделей на рис. 7.5 в разделе 7.3). Эти модели и связанные с ними априорные вероятности, $P(\text{frontier})$, показаны на рис. 12.6. Итак, мы получаем следующие значения:

$$P(P_{1,3} | \text{known}, b) = \alpha' \langle 0,2(0,04 + 0,16 + 0,16), 0,8(0,04 + 0,16) \rangle \approx \langle 0,31, 0,69 \rangle.$$

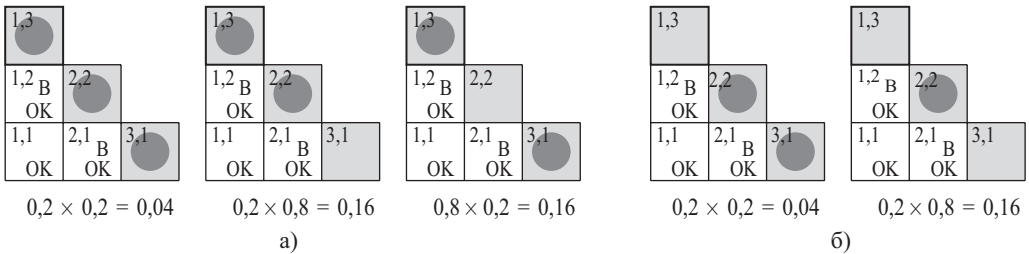


Рис. 12.6. Согласованные модели для переменных множества *Frontier* $P_{2,2}$ и $P_{3,1}$, показывающие значение $P(\text{frontier})$ для каждой модели. **а)** Три модели с $P_{1,3} = \text{true}$, где показаны две или три ямы. **б)** Две модели с $P_{1,3} = \text{false}$, где показаны одна или две ямы

Таким образом, квадрат [1,3] (и квадрат [3,1] по симметрии) содержит яму с вероятностью приблизительно 31%. Аналогичные вычисления, которые читатель вполне может выполнить самостоятельно, показывают, что квадрат [2,2] содержит яму с вероятностью приблизительно 86%. Агент в мире вампуса, определенно, должен избегать квадрата [2,2]! Обратите внимание, что логический агент из главы 7 ничего не знает о том, что выбор квадрата [2,2] может иметь намного худшие последствия, чем выбор любого другого из числа доступных. Логика может лишь сказать нам, что остается неизвестным, есть ли яма в квадрате [2,2], поэтому, чтобы получить больше информации по этому вопросу, необходимо обратиться к вероятностям.

В данном разделе было показано, что даже такие задачи, которые кажутся очень сложными, могут быть точно сформулированы в терминах теории вероятностей и решены с использованием простых алгоритмов. Для получения эффективных решений могут применяться соотношения, определяющие свойства независимости и условной независимости, что позволит упростить необходимые в расчетах операции суммирования. Эти соотношения часто соответствуют нашему интуитивному пониманию того, как следует выполнять декомпозицию задачи. В следующей главе будут разработаны формальные представления для таких соотношений, а также алгоритмы, оперирующие соответствующими представлениями и позволяющие эффективно осуществлять вероятностный логический вывод.

Резюме

В данной главе было предложено использовать теорию вероятности в качестве подходящей основы для рассуждений о неопределенности и дано самое общее представление о возможных способах ее использования.

- Неопределенность возникает как по причине экономии усилий, так и из-за отсутствия знаний. Ее невозможно избежать в сложной, недетерминированной или частично наблюдаемой проблемной среде.
- В оценках **вероятности** выражается неспособность агента прийти к определенному решению в отношении истинности высказывания. Вероятности обобщают степень уверенности агента в отношении свидетельств.
- **Теория принятия решений** объединяет убеждения и намерения агента за счет определения наилучшего действия как максимизирующего ожидаемую **полезность**.
- К основным типам вероятностных высказываний относятся **априорные** или **безусловные вероятности** и **апостериорные** или **условные вероятности** в отношении простых и сложных высказываний.
- Аксиомы вероятностей ограничивают допустимые значения вероятностей логически связанных высказываний. Агент, игнорирующий в своих действиях эти аксиомы, в некоторых обстоятельствах неизбежно будет вести себя нерационально.
- **Полное совместное распределение вероятностей** определяет вероятность каждого полного присваивания значений случайным переменным. Это распределение обычно слишком велико для того, чтобы его можно было создавать или использовать в явной форме, но если оно доступно, то может использоваться для получения ответа на любые запросы простым суммированием значений вероятностей для возможных миров, соответствующих высказываниям запроса.
- **Абсолютная независимость** между подмножествами случайных переменных позволяет выполнить декомпозицию полного совместного распределения на меньшие совместные распределения, что значительно уменьшает его сложность.
- **Правило Байеса** позволяет вычислять неизвестные вероятности из известных условных вероятностей, обычно в причинном направлении. При наличии многочисленных свидетельств применение правила Байеса приводит к возникновению таких же проблем масштабирования, которые возникают при использовании полного совместного распределения.
- Свойство **условной независимости**, вызванное наличием прямых причинных связей в проблемной области, позволяет провести декомпозицию полного совместного распределения на меньшие условные распределения.

В **наивной байесовской** модели предполагается наличие условной независимости всех переменных действия, если задана одна переменная причины; размеры этой модели увеличиваются линейно в зависимости от количества результатов.

- Агент в мире вампуса может вычислять вероятности ненаблюдаемых объектов мира и использовать их для принятия лучших решений в сравнении с простым логическим агентом. Условная независимость делает эти вычисления легко реализуемыми.

Библиографические и исторические заметки

Теория вероятностей появилась как средство анализа азартных игр. Примерно в 850 году н.э. индийский математик Махавирачарья описал, как подобрать набор ставок, исключающий возможность проигрыша (то, что мы сейчас называем *голландской книгой*). Первый значимый систематический анализ был проведен Джироламо Кардано примерно в 1565 году, но его работы были опубликованы только после его смерти (1663). К этому времени вероятность уже сформировалась как математическая дисциплина благодаря серии достижений, о которых Блез Паскаль сообщал в переписке с Пьером Ферма в 1654 году. Первым опубликованным учебником по теории вероятностей была книга Гюйгенса *De Ratiociniis in Ludo Aleae* ([1106], 1657). Взгляд на “лень и невежество” как источник неопределенности был предложен Джоном Арбетнотом в предисловии к его переводу этой книги Гюйгенса ([65], 1692).

Связь между вероятностью и рассуждением восходит по крайней мере к XIX веку: в 1819 году Пьер Лаплас сказал: “Теория вероятностей — это ни что иное, как здравый смысл, сведенный к расчетам”. В 1850 году Джеймс Максвелл сказал: “Истинная логика для этого мира — исчисление вероятностей, которое учитывает величину вероятности, которая есть или должна быть в сознании разумного человека”.

Долгие годы шли бесконечные дебаты по поводу источника и статуса значений вероятности. Сторонники ► **эмпирического** (частотного) подхода к вероятности настаивали на том, что эти значения могут быть получены только из *экспериментов*: если после тестирования 100 человек будет установлено, что десять из них имеют зубную полость, то можно будет утверждать, что вероятность образования такой полости равна приблизительно 0,1. Для сторонников этой точки зрения утверждение “вероятность образования зубной полости равна 0,1” означает, что значение 0,1 представляет собой долю случаев наличия полости, наблюдаемую в пределе бесконечного числа испытаний. Исходя из любой конечной выборки, можно оценить истинную долю, а также вычислить, насколько точной, скорее всего, является эта оценка.

Сторонники ► **объективистского** подхода полагают, что вероятности являются реальными аспектами Вселенной — склонностью самих объектов вести себя

определенным образом, а не просто описанием степени уверенности наблюдателя. Например, тот факт, что для обычной монеты (без жульнических подделок) вероятность выпадения орла составляет 0,5, является склонностью самой монеты. С этой точки зрения любые измерения эмпириков являются просто попытками наблюдать эти склонности. Большинство физиков согласны с тем, что квантовые явления объективно вероятностны, а неопределенность в макроскопическом масштабе — например, при подбрасывании монет — обычно возникает из-за незнания начальных условий и, по-видимому, не согласуется с представлением о склонности.

Сторонники ► **субъективистского** подхода описывают вероятность как способ охарактеризовать убеждения агента, а не как что-то, имеющее некое внешнее физическое значение. Субъективный **байесовский** подход допускает любое самосогласованное приписывание априорных вероятностей высказываниям, однако затем настаивает на их надлежащем байесовском обновлении по мере поступления свидетельств.

Даже строгая эмпирическая позиция предполагает субъективность из-за проблемы ► **эталонного класса**: пытаясь определить вероятность исхода *конкретного* эксперимента, эмпирик должен отнести его к эталонному классу “похожих” экспериментов с известными частотами исхода. Но как выбрать для него правильный класс? И.Дж. Гуд писал: “Каждое событие в жизни уникально, и каждая вероятность в реальной жизни, которую мы оцениваем на практике, — это событие, которое никогда не происходило ранее” (Гуд [894], 1983).

Например, в случае конкретного пациента эмпирик, желающий оценить вероятность наличия у него зубной полости, должен рассмотреть эталонный класс других пациентов, схожих с ним по некоторым важным аспектам — возрасту, симптомам, особенностям питания — и выяснить, какую часть из них составляли те, у кого имелась зубная полость. Если стоматолог примет во внимание все, что ему известно о пациенте, включая цвет волос, вес с точностью до грамма, девичью фамилию матери и так далее, эталонный класс в конечном счете окажется пустым. Эта ситуация была серьезной проблемой в философии науки.

Паскаль использовал вероятность в таких вычислениях, которые требовали не только ее объективной интерпретации как свойства мира, основанного на симметрии или относительных частотах событий, но и субъективной интерпретации, основанной на оценке степени уверенности. Первая интерпретация обнаруживается в проведенном Паскалем анализе вероятностей в играх с элементами случайности, а последняя — в знаменитых доводах “Спора с Паскалем”, касающихся возможного существования Бога. Однако Паскаль недостаточно четко учитывал различие между этими двумя интерпретациями. Указанное различие было впервые наглядно подчеркнуто Джеймсом Бернулли (1654–1705).

Лейбниц ввел “классическое” понятие вероятности как доли перечислимых, равновероятных случаев, которое использовалось также Бернулли, но было полностью проанализировано Лапласом ([1353], 1816). Это понятие является противоречивым из-за наличия эмпирической (частотной) и субъективной интерпретации. События

могут рассматриваться как равновероятные либо из-за наличия естественной, физической симметрии между ними, либо просто из-за того, что мы не обладаем достаточными знаниями, которые позволили бы считать одно событие более вероятным, чем другое. Подход, предусматривающий использование последних, субъективных соображений, оправдывающих допустимость присваивания равных вероятностей, известен под названием ► **принцип безразличия** [792]. Этот принцип часто приписывают Лапласу ([1353], 1816), но он никогда не использовал это название явно, — впервые это сделал Кейнс ([1220], 1921). Джордж Буль и Джон Венн, оба ссылались на него как на ► **принцип недостаточной причины**.

Споры между сторонниками объективного и субъективного подходов еще более обострились в XX столетии. Колмогоров ([1272], 1963), Р.А. Фишер ([745], 1922) и Ричард фон Мизес ([745], 1922) были сторонниками относительной эмпирической (частотной) интерпретации. Приведенная в работе Карла Поппера [1812] (1959) (впервые опубликована на немецком языке в 1934 году) интерпретация “проявлений закономерностей” позволяет проследить истоки формирования относительных частот вплоть до основополагающих законов физической симметрии. Франк Рамсей ([1848], 1931), Бруно де Финетти ([553], 1937), Р.Т. Кокс ([489], 1946), Леонард Сэведж ([1984], 1954), Ричард Джеффри ([1129], 1983) и И.Т. Джейнс (2003) интерпретировали вероятности как степени уверенности конкретных лиц. Их анализ степени уверенности был тесно связан с полезностями и с поведением, а именно — с готовностью субъекта делать те или иные ставки.

Рудольф Карнап предложил иную интерпретацию вероятности — не как определенной степени уверенности конкретного лица, а как степень уверенности, которую идеализированное рассуждающее лицо *должно* иметь в отношении истинности конкретного высказывания *a*, при заданном конкретном ряде свидетельств *e*. Карнап попытался сделать это понятие степени **подтверждения** математически точным, как логическое отношение между *a* и *e*. В настоящее время считается, что уникальной логики подобного типа не существует; скорее, любая подобная логика опирается на субъективное априорное распределение вероятностей, эффект которого уменьшается по мере сбора большего количества наблюдений.

Изучение этого отношения имело целью создание математической дисциплины, названной **индуктивной логикой** по аналогии с обычной дедуктивной логикой (Карнап [373], 1948; [374], 1950). Карнап не смог в достаточной степени расширить свою индуктивную логику за пределы пропозиционального случая, а Патнем ([1829], 1963) на состязательных аргументах показала, что ей присущи некоторые фундаментальные сложности. В более поздней работе Бакхуса, Гроува, Гальперна и Коллера ([99], 1992) метод Карнапа был расширен на теории первого порядка.

Первая строго аксиоматическая основа для теории вероятностей была предложена Колмогоровым ([1271], 1950) (впервые опубликована в Германии в 1933 году). Позднее Рени ([1873], 1970) дал аксиоматическое представление, использующее в качестве примитивов условные, а не абсолютные вероятности.

В дополнение к аргументам де Финетти в отношении обоснованности аксиом Кокса ([489], 1946) показал, что любая система неопределенных рассуждений, соответствующая его набору допущений, эквивалентна теории вероятностей. Это придало поклонникам вероятности новую уверенность, но остальные так и не были убеждены, возражая против допущения, что уверенность должна быть представлена единственным числом. Гальперн ([952], 1999) проанализировал допущения и указал на некоторые пробелы в первоначальной формулировке Кокса. Позднее Хорн ([1062], 2003) показал, как исправить эти трудности, а Джейнс (2003) привел аналогичный аргумент, который легче воспринимается.

Преподобный Томас Байес (1702–1761) сформулировал правило формирования рассуждений об условных вероятностях, которое позднее было названо в его честь (Байес [148], 1763). Но Байес рассматривал только случай равномерных априорных распределений, тогда как Лаплас независимо от него разработал теорию для общего случая. Байесовские вероятностные рассуждения использовались в приложениях ИИ уже с 1960-х годов, особенно в медицинской диагностике. Они использовались не только для постановки диагноза на основе имеющихся данных, но также для выбора необходимых дополнительных вопросов и тестов с использованием теории значения информации (раздел 16.6), когда имеющиеся доказательства были неубедительны (Горри [909], 1968; Горри и др. [910], 1973). Одна система даже превзошла людей-экспертов в диагностике острых брюшных заболеваний (де Домба и др. [550], 1974). В своей статье Лукас и соавт. ([1462], 2004) дали соответствующий обзор.

Однако эти ранние байесовские системы имели множество недостатков. Поскольку в них отсутствовали какие-либо теоретические модели диагностируемых ими условий, они были чувствительны к нерепрезентативным данным, встречающимся в тех ситуациях, когда были доступны лишь небольшие выборки (де Домба и др. [551], 1981). Еще более фундаментальным недостатком было то, что в этих системах не применялись лаконичные формальные средства (подобные тем, которые будут описаны в главе 13) для представления и использования информации об условной независимости информации. Поэтому успешная эксплуатация этих систем зависела от накопления, хранения и обработки громадных таблиц с вероятностными данными. Из-за этих сложностей в период с середины 1970-х до конца 1980-х годов интерес исследователей в области искусственного интеллекта к вероятностным методам решения задач в условиях неопределенности значительно снизился. Новые разработки в этой области, появившиеся лишь в конце 1980-х годов, будут рассмотрены в следующей главе.

Наивная байесовская модель для совместных распределений широко исследовалась в литературе по распознаванию образов уже с 1950-х годов (Дуда и Харт [659], 1973). Кроме того, такой способ представления использовался, часто непреднамеренно, в области выборки информации, начиная с работы Марона ([1496], 1961). Вероятностные основы этого метода, дополнительно рассматриваемые в упражнении 12.28, были исследованы Робертсоном и Спарком Джонсом

([1897], 1976). Домингос и Паццани ([631], 1997) объяснили причины поразительного успеха наивных байесовских рассуждений даже в тех проблемных областях, в которых они явно нарушали предположения о независимости.

По теории вероятностей есть много хороших вводных учебников, включая книги Бертсекаса и Цициклиса ([202], 2008), Росса ([1918], 2015) и Гринстеда и Снелла ([924], 1997). Де Грот и Шервиш ([593], 2001) выпустили объединенное введение в теорию вероятностей и статистику с байесовской точки зрения, а Уолпол и др. ([2289], 2016) предложили вводный курс для ученых и инженеров. Джейнс (2003) дал очень убедительное изложение байесовского подхода. Биллингсли ([215], 2012) и Венкатеш ([2268], 2012) придерживаются более математических методов изложения, включая обсуждение всех осложнений, связанных с непрерывными переменными, которое мы здесь опустили. Хакинг ([941], 1975) и Хелд ([946], 1990) рассматривают также раннюю историю концепции вероятности, а в статье Бенштейна ([192], 1996) приведен популярный обзор.

Упражнения

- 12.1. Исходя из основных принципов докажите, что $P(ab \wedge a) = 1$.
- 12.2. Воспользовавшись аксиомами вероятности, докажите, что любое распределение вероятностей дискретной случайной переменной должно в сумме составлять 1.
- 12.3. Для каждого из следующих высказываний либо докажите, что оно истинно, либо приведите контрпример.
- Если $P(ab, c) = P(ba, c)$, то $P(ac) = P(bc)$
 - Если $P(ab, c) = P(a)$, то $P(bc) = P(b)$
 - Если $P(ab) = P(a)$, то $P(ab, c) = P(ac)$
- 12.4. Будет ли для агента рационально придерживаться трех убеждений: $P(A) = 0,4$; $P(B) = 0,3$ и $P(A \vee B) = 0,5$? Если это так, то какой диапазон вероятностей будет в этом случае рациональным для агента применительно к $A \wedge B$? Составьте таблицу, подобную приведенной на рис. 12.2, и покажите, подтверждает ли она ваши доводы в отношении рациональности. Затем составьте еще одну версию этой таблицы, в которой $P(A \vee B) = 0,7$. Объясните, почему рационально будет принять именно это значение вероятности, несмотря даже на то, что в данной таблице присутствует один случай, соответствующий проигрышу, и три случая с ничейным результатом. (*Подсказка.* Что агент 1 полагает относительно вероятности каждого из этих четырех случаев, в особенности того, когда имеет место проигрыш?)
- 12.5. Этот вопрос касается свойств возможных миров, определенных в разделе 12.2.2 как присваивание значений всем рассматриваемым случайным переменным. Мы будем работать с высказываниями, которые соответствуют точно одному возможному миру, поскольку они включают значения для всех переменных. В теории вероятностей такие высказывания называют **атомарными событиями**. Например, для булевых переменных X_1, X_2, X_3 высказывание $x_1 \wedge \neg x_2 \wedge \neg x_3$

определяет присваивание значений всех переменных, — на языке логики высказываний можно было бы сказать, что у него есть ровно одна модель.

- а) Для случая n булевых переменных докажите, что любые два различных атомарных события являются взаимоисключающими, т.е. их конъюнкция эквивалентна значению *false*.
- б) Докажите, что дизъюнкция всех возможных атомарных событий логически эквивалентна значению *true*.
- в) Докажите, что любое высказывание логически эквивалентно дизъюнкции атомарных событий, которые влекут за собой его истинность.

12.6. Докажите уравнение (12.5) на основании уравнений (12.2) и (12.3).

12.7. Рассмотрим множество из всех возможных раздач по 5 карт при игре в покер со стандартной колодой в 52 карты, полагая, что раздача проводится честно.

- а) Сколько атомарных событий будет в совместном распределении вероятностей (т.е. сколько существует различных раздач по пять карт)?
- б) Какова вероятность каждого атомарного события?
- в) Какова вероятность получения королевского флеш-стрита? А любого каре из всех возможных?

12.8. Исходя из полного совместного распределения, приведенного на рис. 12.3, рассчитайте следующее.

- а) $P(\textit{toothache})$
- б) $P(\textit{Cavity})$
- в) $P(\textit{Toothache} \mid \textit{cavity})$
- г) $P(\textit{Cavity} \mid \textit{toothache} \vee \textit{catch})$

12.9. Исходя из полного совместного распределения, приведенного на рис. 12.3, рассчитайте следующее.

- а) $P(\textit{toothache})$
- б) $P(\textit{Catch})$
- в) $P(\textit{Cavity} \mid \textit{catch})$
- г) $P(\textit{Cavity} \mid \textit{toothache} \vee \textit{catch})$

12.10. В своем письме от 24 августа 1654 года Паскаль попытался показать, как следует распределять денежные ставки, когда азартная игра должна закончиться преждевременно. Представьте себе игру, в которой каждый ход состоит из броска игральной кости. Игрок E получает очко, если выпадает четное число, а игрок O получает очко, если выпавшее число нечетное. Первый из игроков, набравший 7 очков, выигрывает банк. Предположим, что игра прерывается, когда игрок E ведет со счетом 4:2. Как в этом случае справедливо разделить деньги в банке между игроками? Какова будет общая формула? (Ферма и Паскаль сделали несколько ошибок, прежде чем решить проблему, но вы должны быть в состоянии сделать это правильно с первого раза.)

12.11. Решив использовать теорию вероятностей на практике, мы выбрали игровой автомат с тремя независимыми колесами, каждое из которых с равной вероятностью может отображать один из четырех символов: золотой слиток, колокольчик, лимон или вишню. Игровой автомат имеет следующую схему выплат для

ставки в одну монету (здесь “?” означает, что не имеет значения, что отображается на данном колесе):

- золотой слиток/золотой слиток/золотой слиток — выплата 20 монет
- колокольчик/колокольчик/колокольчик — выплата 15 монет
- лимон/лимон/лимон — выплата 5 монет
- вишня/вишня/вишня — выплата 3 монет
- вишня/вишня/? — выплата 2 монет
- вишня/?/? — выплата 1 монеты

- а) Рассчитайте ожидаемый процент “окупаемости” автомата. Другими словами, какова ожидаемая выдача монет для каждой монеты, заплаченной за игру?
- б) Вычислите вероятность того, что игра на этом игровом автомате приведет к выигрышу.
- в) Оцените среднее и медианное количество игр, которое, как можно ожидать, будет сыграно до проигрыша всей имеющейся суммы, если начать с 10 монет. Можете запустить симуляцию, чтобы просто оценить эти значения, вместо того чтобы пытаться вычислить точный ответ.

12.12. Решив использовать теорию вероятностей на практике, мы выбрали игровой автомат с тремя независимыми колесами, каждое из которых с равной вероятностью может отображать один из четырех символов: золотой слиток, колокольчик, лимон или вишню. Игровой автомат имеет следующую схему выплат для ставки в одну монету (здесь “?” означает, что не имеет значения, что отображается на данном колесе):

- золотой слиток/золотой слиток/золотой слиток — выплата 21 монеты
- колокольчик/колокольчик/колокольчик — выплата 16 монет
- лимон/лимон/лимон — выплата 5 монет
- вишня/вишня/вишня — выплата 3 монет
- вишня/вишня/? — выплата 2 монет
- вишня/?/? — выплата 1 монеты

- а) Рассчитайте ожидаемый процент “окупаемости” автомата. Другими словами, какова ожидаемая выдача монет для каждой монеты, заплаченной за игру?
- б) Вычислите вероятность того, что игра на этом игровом автомате приведет к выигрышу.
- в) Оцените среднее и медианное количество игр, которое, как можно ожидать, будет сыграно до проигрыша всей имеющейся суммы, если начать с 8 монет. Можете запустить симуляцию, чтобы просто оценить эти значения, вместо того чтобы пытаться вычислить точный ответ.

12.13. Необходимо передать n -битное сообщение агенту-получателю. В процессе передачи биты в сообщении независимо повреждаются (меняют значение на противоположное) с вероятностью ϵ для каждого. Используя дополнительный бит четности, отправляемый вместе с исходной информацией, получатель сможет восстановить сообщение, если повреждено не более одного бита во всем сообщении (включая бит четности). Предположим, нужно убедиться, что правильное сообщение будет получено с вероятностью не менее $1 - \delta$. Каким в этом

случае будет максимально допустимое значение n ? Рассчитайте это значение для случая $\epsilon = 0,001$ и $\delta = 0,01$.

- 12.14.** Необходимо передать n -битное сообщение агенту-получателю. В процессе передачи биты в сообщении независимо повреждаются (меняют значение на противоположное) с вероятностью ϵ для каждого. Используя дополнительный бит четности, отправляемый вместе с исходной информацией, получатель сможет восстановить сообщение, если повреждено не более одного бита во всем сообщении (включая бит четности). Предположим, нужно убедиться, что правильное сообщение будет получено с вероятностью не менее $1 - \delta$. Каким в этом случае будет максимально допустимое значение n ? Рассчитайте это значение для случая $\epsilon = 0,002$ и $\delta = 0,01$.
- 12.15.** Покажите, что три формы описания свойства независимости, приведенные в уравнении (12.11), являются эквивалентными.
- 12.16.** Рассмотрим два медицинских теста на некоторый вирус, А и В. Тест А обладает эффективностью 95% при обнаружении вируса, когда он действительно присутствует, но дает 10% ложных положительных результатов (указывает, что вирус присутствует, когда на самом деле его нет). Тест В обладает эффективностью 90% при обнаружении вируса и дает 5% ложных положительных результатов. В этих двух тестах используются независимые методы идентификации вируса. Вирус присутствует у 1% всех людей. Пусть каждого человека проверяют на наличие вируса, используя только один из тестов, и для переносчиков вируса он дает положительный результат. Для какого из двух тестов получение положительного результата будет более показательным, если человек действительно является переносчиком вируса? Обоснуйте свой ответ математически.
- 12.17.** Предположим, дана монета, которая падает орлом вверх с вероятностью x и решкой вверх — с вероятностью $1 - x$. Являются ли результаты последовательных бросков монеты независимыми друг от друга, если вы *знаете* значение x ? Являются ли результаты последовательных бросков монеты независимыми друг от друга, если вы *не знаете* значение x ? Обоснуйте свой ответ.
- 12.18.** После ежегодного медицинского осмотра пациента у врача есть плохая новость и хорошая новость. Плохая новость состоит в том, что проверка на наличие серьезного заболевания оказалась положительной и что точность результатов этой проверки составляет 99% (т.е. вероятность получения положительного результата проверки, если пациент имеет это заболевание, равна 0,99, и такова же вероятность получения отрицательных результатов проверки, если пациент не имеет этого заболевания). Хорошая новость состоит в том, что это заболевание — редкое и поражает только одного из 10 тысяч людей того возраста, в котором находится пациент. Почему новость, что это заболевание редкое, следует считать хорошей? Каковы шансы на то, что пациент действительно имеет данное заболевание?
- 12.19.** После ежегодного медицинского осмотра пациента у врача есть плохая новость и хорошая новость. Плохая новость состоит в том, что проверка на наличие серьезного заболевания оказалась положительной и что точность результатов этой проверки

составляет 99% (т.е. вероятность получения положительного результата проверки, если пациент имеет это заболевание, равна 0,99, и такова же вероятность получения отрицательных результатов проверки, если пациент не имеет этого заболевания). Хорошая новость состоит в том, что это заболевание — редкое и поражает только одного из 100 тысяч людей того возраста, в котором находится пациент. Почему новость, что это заболевание редкое, следует считать хорошей? Каковы шансы на то, что пациент действительно имеет данное заболевание?

12.20. Довольно часто полезно рассмотреть результаты некоторых конкретных высказываний в контексте некоторого общего фоновое свидетельства, которое остается неизменным, а не действовать в условиях полного отсутствия информации. В приведенных ниже вопросах предлагается доказать более общие версии правила умножения вероятностей и правила Байеса применительно к некоторому фоновому свидетельству e .

а) Докажите версию общего правила умножения вероятностей для условных вероятностей:

$$\mathbf{P}(X, Ye) = \mathbf{P}(XY, e)\mathbf{P}(Ye).$$

б) Докажите версию правила Байеса с условными вероятностями из уравнения (12.13).

12.21. Покажите, что утверждение условной независимости

$$\mathbf{P}(X, Y | Z) = \mathbf{P}(X | Z)\mathbf{P}(Y | Z)$$

эквивалентно любому из следующих утверждений:

$$\mathbf{P}(X | Y, Z) = \mathbf{P}(X | Z) \quad \text{и} \quad \mathbf{P}(Y | X, Z) = \mathbf{P}(Y | Z).$$

12.22. Предположим, вам вручили мешок, содержащий n подлинных монет, и сообщили, что $n - 1$ из этих монет являются нормальными, т.е. такими, что с одной стороны у них орел, с другой — решка, а одна монета — фальшивая: на обеих ее сторонах изображен орел.

а) Допустим, вы открыли мешок, случайным образом выбрали монету и подбросили ее, в результате чего выпал орел. Какова (условная) вероятность того, что выбранная вами монета является фальшивой?

б) Теперь предположим, что вы продолжаете подбрасывать эту монету в общей сложности k раз после того, как она была выбрана, и наблюдаете k выпадений орла. Какова теперь условная вероятность того, что вы выбрали фальшивую монету?

в) Наконец, предположим, что вы хотите принять решение, является ли выбранная монета фальшивой, подбросив ее k раз. Процедура принятия решения возвращает *fake* (фальшивая), если все k бросков приводят к выпадению орла, а в противном случае она возвращает значение *normal* (нормальная). Какова (безусловная) вероятность того, что эта процедура сделает ошибку?

12.23. В этом упражнении требуется вычислить коэффициент нормализации для примера с заболеванием менингитом. Вначале выберите подходящее значение для $P(s | \neg m)$ и примените его для вычисления ненормализованных значений $P(m | s)$ и $P(\neg m | s)$ (т.е. игнорируя терм $P(s)$ в выражении правила Байеса). Теперь нормализуйте эти значения таким образом, чтобы они в сумме составляли 1.

- 12.24.** В этом упражнении исследуется то, как соотношения, касающиеся условной независимости, влияют на количество информации, требуемой для вероятностных вычислений.
- а) Предположим, что необходимо рассчитать значение $P(h | e_1, e_2)$, а информация об условной независимости отсутствует. Какие из следующих множеств чисел являются достаточными для такого вычисления?
1. $P(E_1, E_2)$, $P(H)$, $P(E_1 | H)$, $P(E_2 | H)$
 2. $P(E_1, E_2)$, $P(H)$, $P(E_1, E_2 | H)$
 3. $P(H)$, $P(E_1 | H)$, $P(E_2 | H)$
- б) Предположим, известно, что $P(E_1 | H, E_2) = P(E_1 | H)$ для всех значений H, E_1, E_2 . Какое из этих трех множеств значений теперь будет достаточным?
- 12.25.** Пусть X, Y, Z — случайные булевы переменные. Обозначьте восемь элементов совместного распределения $\mathbf{P}(X, Y, Z)$ буквами алфавита от a до h . Выразите утверждение, что X и Y являются условно независимыми при заданном Z в виде множества уравнений, связывающих элементы от a до h . Сколько среди них *неизбыточных* уравнений?
- 12.26.** Предположим, что вы — свидетель ночного наезда на пешехода в Афинах с участием такси, которое скрылось с места происшествия. Все такси в Афинах покрашены в синий или зеленый цвет. Вы поклялись под присягой, что такси было синим, при этом результаты широких экспериментов показывают, что в условиях плохого освещения надежность распознавания синего и зеленого цветов составляет 75%.
- а) Возможно ли рассчитать наиболее вероятный цвет этого такси? (*Подсказка.* Тщательно проведите различие между высказыванием, что такси — синего цвета, и высказыванием, что оно показалось вам синим.)
- а) Что изменится после получения информации о том, что в Афинах 9 из 10 такси зеленого цвета?
- 12.27.** Запишите общий алгоритм получения ответов на запросы в форме $\mathbf{P}(Cause | e)$, используя наивное байесовское распределение. Исходите из предположения, что свидетельство e может присваивать значения любому подмножеству переменных результата.
- 12.28.** Категоризацией текста называется задача присваивания данному конкретному документу одной из фиксированного множества категорий на основе анализа его текста. Для решения этой задачи часто используются наивные байесовские модели, в которых переменной запроса является категория документа, а в качестве переменных “результата” рассматривается наличие или отсутствие каждого слова в “языке” категории. Основное предположение состоит в том, что слова в документах встречаются независимо друг от друга, а их частоты определяются категорией документа.
- а) Дайте точное объяснение, как можно сформировать такую модель, получив в качестве “обучающих данных” множество документов, уже распределенных по категориям.

- б) Дайте точное объяснение, как следует определять категорию нового документа.
- в) Является ли указанное предположение о независимости обоснованным? Обсудите этот вопрос.

- 12.29.** В проведенном в этой главе анализе мира вампуса использовался тот факт, что каждый квадрат содержит яму с вероятностью 0,2, независимо от содержимого других квадратов. Вместо этого примем предположение, что точно $N/5$ ям равномерно разбросаны случайным образом среди N квадратов, отличных от [1,1]. Останутся ли переменные $P_{i,j}$ и $P_{k,l}$ все еще независимыми? Каково теперь совместное распределение $\mathbf{P}(P_{1,1}, \dots, P_{4,4})$? Заново выполните вычисление вероятностей наличия ям в квадратах [1,3] и [2,2].
- 12.30.** Повторите расчет вероятности наличия ям в квадратах [1,3] и [2,2], полагая, что каждый квадрат содержит яму с вероятностью 0,01, независимо от других квадратов. Что в этом случае можно будет сказать об относительной эффективности логического и вероятностного агента?
- 12.31.** Реализуйте гибридного вероятностного агента для мира вампуса, основываясь на гибридном агенте, представленном на рис. 7.20, и процедуре вероятностного вывода, рассмотренной в этой главе.