

ables for squares other than the *Known* squares and the query square [1,3]. Then, by Equation (13.9), we have

$$\mathbf{P}(P_{1,3} \mid \text{known}, b) = \alpha \sum_{\text{unknown}} \mathbf{P}(P_{1,3}, \text{unknown}, \text{known}, b) .$$

The full joint probabilities have already been specified, so we are done—that is, unless we care about computation. There are 12 unknown squares; hence the summation contains  $2^{12} = 4096$  terms. In general, the summation grows exponentially with the number of squares.

Surely, one might ask, aren't the other squares irrelevant? How could [4,4] affect whether [1,3] has a pit? Indeed, this intuition is correct. Let *Frontier* be the pit variables (other than the query variable) that are adjacent to visited squares, in this case just [2,2] and [3,1]. Also, let *Other* be the pit variables for the other unknown squares; in this case, there are 10 other squares, as shown in Figure 13.5(b). The key insight is that the observed breezes are *conditionally independent* of the other variables, given the known, frontier, and query variables. To use the insight, we manipulate the query formula into a form in which the breezes are conditioned on all the other variables, and then we apply conditional independence:

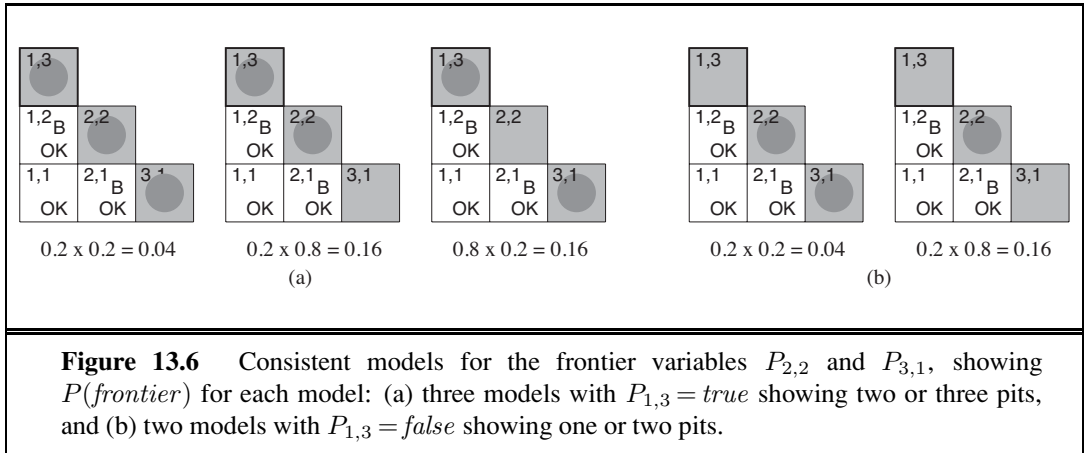
$$\begin{aligned} \mathbf{P}(P_{1,3} \mid \text{known}, b) &= \alpha \sum_{\text{unknown}} \mathbf{P}(P_{1,3}, \text{known}, b, \text{unknown}) \quad (\text{by Equation (13.9)}) \\ &= \alpha \sum_{\text{unknown}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{unknown}) \mathbf{P}(P_{1,3}, \text{known}, \text{unknown}) \\ &\quad (\text{by the product rule}) \\ &= \alpha \sum_{\text{frontier}} \sum_{\text{other}} \mathbf{P}(b \mid \text{known}, P_{1,3}, \text{frontier}, \text{other}) \mathbf{P}(P_{1,3}, \text{known}, \text{frontier}, \text{other}) \\ &= \alpha \sum_{\text{frontier}} \sum_{\text{other}} \mathbf{P}(b \mid \text{known}, P_{1,3}, \text{frontier}) \mathbf{P}(P_{1,3}, \text{known}, \text{frontier}, \text{other}) , \end{aligned}$$

where the final step uses conditional independence: *b* is independent of *other* given *known*, *P*<sub>1,3</sub>, and *frontier*. Now, the first term in this expression does not depend on the *Other* variables, so we can move the summation inward:

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

By independence, as in Equation (13.20), the prior term can be factored, and then the terms can be reordered:

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



where the last step folds  $P(\text{known})$  into the normalizing constant and uses the fact that  $\sum_{\text{other}} P(\text{other})$  equals 1.

Now, there are just four terms in the summation over the frontier variables  $P_{2,2}$  and  $P_{3,1}$ . The use of independence and conditional independence has completely eliminated the other squares from consideration.

Notice that the expression  $\mathbf{P}(b \mid \text{known}, P_{1,3}, \text{frontier})$  is 1 when the frontier is consistent with the breeze observations, and 0 otherwise. Thus, for each value of  $P_{1,3}$ , we sum over the *logical models* for the frontier variables that are consistent with the known facts. (Compare with the enumeration over models in Figure 7.5 on page 241.) The models and their associated prior probabilities— $P(\text{frontier})$ —are shown in Figure 13.6. We have

$$\mathbf{P}(P_{1,3} \mid \text{known}, b) = \alpha' \langle 0.2(0.04 + 0.16 + 0.16), 0.8(0.04 + 0.16) \rangle \approx \langle 0.31, 0.69 \rangle.$$

That is, [1,3] (and [3,1] by symmetry) contains a pit with roughly 31% probability. A similar calculation, which the reader might wish to perform, shows that [2,2] contains a pit with roughly 86% probability. The wumpus agent should definitely avoid [2,2]! Note that our logical agent from Chapter 7 did not know that [2,2] was worse than the other squares. Logic can tell us that it is unknown whether there is a pit in [2, 2], but we need probability to tell us how likely it is.

What this section has shown is that even seemingly complicated problems can be formulated precisely in probability theory and solved with simple algorithms. To get *efficient* solutions, independence and conditional independence relationships can be used to simplify the summations required. These relationships often correspond to our natural understanding of how the problem should be decomposed. In the next chapter, we develop formal representations for such relationships as well as algorithms that operate on those representations to perform probabilistic inference efficiently.

---

## 13.7 SUMMARY

---

This chapter has suggested probability theory as a suitable foundation for uncertain reasoning and provided a gentle introduction to its use.

- Uncertainty arises because of both laziness and ignorance. It is inescapable in complex, nondeterministic, or partially observable environments.
- Probabilities express the agent's inability to reach a definite decision regarding the truth of a sentence. Probabilities summarize the agent's beliefs relative to the evidence.
- Decision theory combines the agent's beliefs and desires, defining the best action as the one that maximizes expected utility.
- Basic probability statements include **prior probabilities** and **conditional probabilities** over simple and complex propositions.
- The axioms of probability constrain the possible assignments of probabilities to propositions. An agent that violates the axioms must behave irrationally in some cases.
- The **full joint probability distribution** specifies the probability of each complete assignment of values to random variables. It is usually too large to create or use in its explicit form, but when it is available it can be used to answer queries simply by adding up entries for the possible worlds corresponding to the query propositions.
- **Absolute independence** between subsets of random variables allows the full joint distribution to be factored into smaller joint distributions, greatly reducing its complexity. Absolute independence seldom occurs in practice.
- **Bayes' rule** allows unknown probabilities to be computed from known conditional probabilities, usually in the causal direction. Applying Bayes' rule with many pieces of evidence runs into the same scaling problems as does the full joint distribution.
- **Conditional independence** brought about by direct causal relationships in the domain might allow the full joint distribution to be factored into smaller, conditional distributions. The **naive Bayes** model assumes the conditional independence of all effect variables, given a single cause variable, and grows linearly with the number of effects.
- A wumpus-world agent can calculate probabilities for unobserved aspects of the world, thereby improving on the decisions of a purely logical agent. Conditional independence makes these calculations tractable.

---

## BIBLIOGRAPHICAL AND HISTORICAL NOTES

Probability theory was invented as a way of analyzing games of chance. In about 850 A.D. the Indian mathematician Mahaviracarya described how to arrange a set of bets that can't lose (what we now call a Dutch book). In Europe, the first significant systematic analyses were produced by Girolamo Cardano around 1565, although publication was posthumous (1663). By that time, probability had been established as a mathematical discipline due to a series of

results established in a famous correspondence between Blaise Pascal and Pierre de Fermat in 1654. As with probability itself, the results were initially motivated by gambling problems (see Exercise 13.9). The first published textbook on probability was *De Ratiociniis in Ludo Aleae* (Huygens, 1657). The “laziness and ignorance” view of uncertainty was described by John Arbuthnot in the preface of his translation of Huygens (Arbuthnot, 1692): “It is impossible for a Die, with such determin’d force and direction, not to fall on such determin’d side, only I don’t know the force and direction which makes it fall on such determin’d side, and therefore I call it Chance, which is nothing but the want of art...”

Laplace (1816) gave an exceptionally accurate and modern overview of probability; he was the first to use the example “take two urns, A and B, the first containing four white and two black balls, ...” The Rev. Thomas Bayes (1702–1761) introduced the rule for reasoning about conditional probabilities that was named after him (Bayes, 1763). Bayes only considered the case of uniform priors; it was Laplace who independently developed the general case. Kolmogorov (1950, first published in German in 1933) presented probability theory in a rigorously axiomatic framework for the first time. Rényi (1970) later gave an axiomatic presentation that took conditional probability, rather than absolute probability, as primitive.

Pascal used probability in ways that required both the objective interpretation, as a property of the world based on symmetry or relative frequency, and the subjective interpretation, based on degree of belief—the former in his analyses of probabilities in games of chance, the latter in the famous “Pascal’s wager” argument about the possible existence of God. However, Pascal did not clearly realize the distinction between these two interpretations. The distinction was first drawn clearly by James Bernoulli (1654–1705).

Leibniz introduced the “classical” notion of probability as a proportion of enumerated, equally probable cases, which was also used by Bernoulli, although it was brought to prominence by Laplace (1749–1827). This notion is ambiguous between the frequency interpretation and the subjective interpretation. The cases can be thought to be equally probable either because of a natural, physical symmetry between them, or simply because we do not have any knowledge that would lead us to consider one more probable than another. The use of this latter, subjective consideration to justify assigning equal probabilities is known as the **principle of indifference**. The principle is often attributed to Laplace, but he never isolated the principle explicitly. George Boole and John Venn both referred to it as the **principle of insufficient reason**; the modern name is due to Keynes (1921).

The debate between objectivists and subjectivists became sharper in the 20th century. Kolmogorov (1963), R. A. Fisher (1922), and Richard von Mises (1928) were advocates of the relative frequency interpretation. Karl Popper’s (1959, first published in German in 1934) “propensity” interpretation traces relative frequencies to an underlying physical symmetry. Frank Ramsey (1931), Bruno de Finetti (1937), R. T. Cox (1946), Leonard Savage (1954), Richard Jeffrey (1983), and E. T. Jaynes (2003) interpreted probabilities as the degrees of belief of specific individuals. Their analyses of degree of belief were closely tied to utilities and to behavior—specifically, to the willingness to place bets. Rudolf Carnap, following Leibniz and Laplace, offered a different kind of subjective interpretation of probability—not as any actual individual’s degree of belief, but as the degree of belief that an idealized individual *should* have in a particular proposition *a*, given a particular body of evidence *e*.

PRINCIPLE OF  
INDIFFERENCE

PRINCIPLE OF  
INSUFFICIENT  
REASON

CONFIRMATION

INDUCTIVE LOGIC

Carnap attempted to go further than Leibniz or Laplace by making this notion of degree of **confirmation** mathematically precise, as a logical relation between  $a$  and  $e$ . The study of this relation was intended to constitute a mathematical discipline called **inductive logic**, analogous to ordinary deductive logic (Carnap, 1948, 1950). Carnap was not able to extend his inductive logic much beyond the propositional case, and Putnam (1963) showed by adversarial arguments that some fundamental difficulties would prevent a strict extension to languages capable of expressing arithmetic.

Cox's theorem (1946) shows that any system for uncertain reasoning that meets his set of assumptions is equivalent to probability theory. This gave renewed confidence to those who already favored probability, but others were not convinced, pointing to the assumptions (primarily that belief must be represented by a single number, and thus the belief in  $\neg p$  must be a function of the belief in  $p$ ). Halpern (1999) describes the assumptions and shows some gaps in Cox's original formulation. Horn (2003) shows how to patch up the difficulties. Jaynes (2003) has a similar argument that is easier to read.

The question of reference classes is closely tied to the attempt to find an inductive logic. The approach of choosing the "most specific" reference class of sufficient size was formally proposed by Reichenbach (1949). Various attempts have been made, notably by Henry Kyburg (1977, 1983), to formulate more sophisticated policies in order to avoid some obvious fallacies that arise with Reichenbach's rule, but such approaches remain somewhat *ad hoc*. More recent work by Bacchus, Grove, Halpern, and Koller (1992) extends Carnap's methods to first-order theories, thereby avoiding many of the difficulties associated with the straightforward reference-class method. Kyburg and Teng (2006) contrast probabilistic inference with nonmonotonic logic.

Bayesian probabilistic reasoning has been used in AI since the 1960s, especially in medical diagnosis. It was used not only to make a diagnosis from available evidence, but also to select further questions and tests by using the theory of information value (Section 16.6) when available evidence was inconclusive (Gorry, 1968; Gorry *et al.*, 1973). One system outperformed human experts in the diagnosis of acute abdominal illnesses (de Dombal *et al.*, 1974). Lucas *et al.* (2004) gives an overview. These early Bayesian systems suffered from a number of problems, however. Because they lacked any theoretical model of the conditions they were diagnosing, they were vulnerable to unrepresentative data occurring in situations for which only a small sample was available (de Dombal *et al.*, 1981). Even more fundamentally, because they lacked a concise formalism (such as the one to be described in Chapter 14) for representing and using conditional independence information, they depended on the acquisition, storage, and processing of enormous tables of probabilistic data. Because of these difficulties, probabilistic methods for coping with uncertainty fell out of favor in AI from the 1970s to the mid-1980s. Developments since the late 1980s are described in the next chapter.

The naive Bayes model for joint distributions has been studied extensively in the pattern recognition literature since the 1950s (Duda and Hart, 1973). It has also been used, often unwittingly, in information retrieval, beginning with the work of Maron (1961). The probabilistic foundations of this technique, described further in Exercise 13.22, were elucidated by Robertson and Sparck Jones (1976). Domingos and Pazzani (1997) provide an explanation

for the surprising success of naive Bayesian reasoning even in domains where the independence assumptions are clearly violated.

There are many good introductory textbooks on probability theory, including those by Bertsekas and Tsitsiklis (2008) and Grinstead and Snell (1997). DeGroot and Schervish (2001) offer a combined introduction to probability and statistics from a Bayesian standpoint. Richard Hamming's (1991) textbook gives a mathematically sophisticated introduction to probability theory from the standpoint of a propensity interpretation based on physical symmetry. Hacking (1975) and Hald (1990) cover the early history of the concept of probability. Bernstein (1996) gives an entertaining popular account of the story of risk.

---

## EXERCISES

**13.1** Show from first principles that  $P(a | b \wedge a) = 1$ .

**13.2** Using the axioms of probability, prove that any probability distribution on a discrete random variable must sum to 1.

**13.3** For each of the following statements, either prove it is true or give a counterexample.

- a. If  $P(a | b, c) = P(b | a, c)$ , then  $P(a | c) = P(b | c)$
- b. If  $P(a | b, c) = P(a)$ , then  $P(b | c) = P(b)$
- c. If  $P(a | b) = P(a)$ , then  $P(a | b, c) = P(a | c)$

**13.4** Would it be rational for an agent to hold the three beliefs  $P(A) = 0.4$ ,  $P(B) = 0.3$ , and  $P(A \vee B) = 0.5$ ? If so, what range of probabilities would be rational for the agent to hold for  $A \wedge B$ ? Make up a table like the one in Figure 13.2, and show how it supports your argument about rationality. Then draw another version of the table where  $P(A \vee B) = 0.7$ . Explain why it is rational to have this probability, even though the table shows one case that is a loss and three that just break even. (*Hint*: what is Agent 1 committed to about the probability of each of the four cases, especially the case that is a loss?)

**13.5** This question deals with the properties of possible worlds, defined on page 488 as assignments to all random variables. We will work with propositions that correspond to exactly one possible world because they pin down the assignments of all the variables. In probability theory, such propositions are called **atomic events**. For example, with Boolean variables  $X_1, X_2, X_3$ , the proposition  $x_1 \wedge \neg x_2 \wedge \neg x_3$  fixes the assignment of the variables; in the language of propositional logic, we would say it has exactly one model.

- a. Prove, for the case of  $n$  Boolean variables, that any two distinct atomic events are mutually exclusive; that is, their conjunction is equivalent to *false*.
- b. Prove that the disjunction of all possible atomic events is logically equivalent to *true*.
- c. Prove that any proposition is logically equivalent to the disjunction of the atomic events that entail its truth.

**13.6** Prove Equation (13.4) from Equations (13.1) and (13.2).

**13.7** Consider the set of all possible five-card poker hands dealt fairly from a standard deck of fifty-two cards.

- How many atomic events are there in the joint probability distribution (i.e., how many five-card hands are there)?
- What is the probability of each atomic event?
- What is the probability of being dealt a royal straight flush? Four of a kind?

**13.8** Given the full joint distribution shown in Figure 13.3, calculate the following:

- $P(\textit{toothache})$ .
- $P(\textit{Cavity})$ .
- $P(\textit{Toothache} \mid \textit{cavity})$ .
- $P(\textit{Cavity} \mid \textit{toothache} \vee \textit{catch})$ .

**13.9** In his letter of August 24, 1654, Pascal was trying to show how a pot of money should be allocated when a gambling game must end prematurely. Imagine a game where each turn consists of the roll of a die, player *E* gets a point when the die is even, and player *O* gets a point when the die is odd. The first player to get 7 points wins the pot. Suppose the game is interrupted with *E* leading 4–2. How should the money be fairly split in this case? What is the general formula? (Fermat and Pascal made several errors before solving the problem, but you should be able to get it right the first time.)

**13.10** Deciding to put probability theory to good use, we encounter a slot machine with three independent wheels, each producing one of the four symbols BAR, BELL, LEMON, or CHERRY with equal probability. The slot machine has the following payout scheme for a bet of 1 coin (where “?” denotes that we don’t care what comes up for that wheel):

BAR/BAR/BAR pays 20 coins  
 BELL/BELL/BELL pays 15 coins  
 LEMON/LEMON/LEMON pays 5 coins  
 CHERRY/CHERRY/CHERRY pays 3 coins  
 CHERRY/CHERRY/? pays 2 coins  
 CHERRY/?/? pays 1 coin

- Compute the expected “payback” percentage of the machine. In other words, for each coin played, what is the expected coin return?
- Compute the probability that playing the slot machine once will result in a win.
- Estimate the mean and median number of plays you can expect to make until you go broke, if you start with 10 coins. You can run a simulation to estimate this, rather than trying to compute an exact answer.

**13.11** We wish to transmit an  $n$ -bit message to a receiving agent. The bits in the message are independently corrupted (flipped) during transmission with  $\epsilon$  probability each. With an extra parity bit sent along with the original information, a message can be corrected by the receiver

if at most one bit in the entire message (including the parity bit) has been corrupted. Suppose we want to ensure that the correct message is received with probability at least  $1 - \delta$ . What is the maximum feasible value of  $n$ ? Calculate this value for the case  $\epsilon = 0.001$ ,  $\delta = 0.01$ .

**13.12** Show that the three forms of independence in Equation (13.11) are equivalent.

**13.13** Consider two medical tests, A and B, for a virus. Test A is 95% effective at recognizing the virus when it is present, but has a 10% false positive rate (indicating that the virus is present, when it is not). Test B is 90% effective at recognizing the virus, but has a 5% false positive rate. The two tests use independent methods of identifying the virus. The virus is carried by 1% of all people. Say that a person is tested for the virus using only one of the tests, and that test comes back positive for carrying the virus. Which test returning positive is more indicative of someone really carrying the virus? Justify your answer mathematically.

**13.14** Suppose you are given a coin that lands *heads* with probability  $x$  and *tails* with probability  $1 - x$ . Are the outcomes of successive flips of the coin independent of each other given that you know the value of  $x$ ? Are the outcomes of successive flips of the coin independent of each other if you do *not* know the value of  $x$ ? Justify your answer.

**13.15** After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease and that the test is 99% accurate (i.e., the probability of testing positive when you do have the disease is 0.99, as is the probability of testing negative when you don't have the disease). The good news is that this is a rare disease, striking only 1 in 10,000 people of your age. Why is it good news that the disease is rare? What are the chances that you actually have the disease?

**13.16** It is quite often useful to consider the effect of some specific propositions in the context of some general background evidence that remains fixed, rather than in the complete absence of information. The following questions ask you to prove more general versions of the product rule and Bayes' rule, with respect to some background evidence  $\mathbf{e}$ :

- a. Prove the conditionalized version of the general product rule:

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

- b. Prove the conditionalized version of Bayes' rule in Equation (13.13).

**13.17** Show that the statement of conditional independence

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

is equivalent to each of the statements

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

**13.18** Suppose you are given a bag containing  $n$  unbiased coins. You are told that  $n - 1$  of these coins are normal, with heads on one side and tails on the other, whereas one coin is a fake, with heads on both sides.

- a. Suppose you reach into the bag, pick out a coin at random, flip it, and get a head. What is the (conditional) probability that the coin you chose is the fake coin?



- b. Suppose you continue flipping the coin for a total of  $k$  times after picking it and see  $k$  heads. Now what is the conditional probability that you picked the fake coin?
- c. Suppose you wanted to decide whether the chosen coin was fake by flipping it  $k$  times. The decision procedure returns *fake* if all  $k$  flips come up heads; otherwise it returns *normal*. What is the (unconditional) probability that this procedure makes an error?

**13.19** In this exercise, you will complete the normalization calculation for the meningitis example. First, make up a suitable value for  $P(s | \neg m)$ , and use it to calculate unnormalized values for  $P(m | s)$  and  $P(\neg m | s)$  (i.e., ignoring the  $P(s)$  term in the Bayes' rule expression, Equation (13.14)). Now normalize these values so that they add to 1.

**13.20** Let  $X, Y, Z$  be Boolean random variables. Label the eight entries in the joint distribution  $\mathbf{P}(X, Y, Z)$  as  $a$  through  $h$ . Express the statement that  $X$  and  $Y$  are conditionally independent given  $Z$ , as a set of equations relating  $a$  through  $h$ . How many *nonredundant* equations are there?

**13.21** (Adapted from Pearl (1988).) Suppose you are a witness to a nighttime hit-and-run accident involving a taxi in Athens. All taxis in Athens are blue or green. You swear, under oath, that the taxi was blue. Extensive testing shows that, under the dim lighting conditions, discrimination between blue and green is 75% reliable.

- a. Is it possible to calculate the most likely color for the taxi? (*Hint*: distinguish carefully between the proposition that the taxi *is* blue and the proposition that it *appears* blue.)
- b. What if you know that 9 out of 10 Athenian taxis are green?

**13.22** Text categorization is the task of assigning a given document to one of a fixed set of categories on the basis of the text it contains. Naive Bayes models are often used for this task. In these models, the query variable is the document category, and the “effect” variables are the presence or absence of each word in the language; the assumption is that words occur independently in documents, with frequencies determined by the document category.

- a. Explain precisely how such a model can be constructed, given as “training data” a set of documents that have been assigned to categories.
- b. Explain precisely how to categorize a new document.
- c. Is the conditional independence assumption reasonable? Discuss.

**13.23** In our analysis of the wumpus world, we used the fact that each square contains a pit with probability 0.2, independently of the contents of the other squares. Suppose instead that exactly  $N/5$  pits are scattered at random among the  $N$  squares other than  $[1,1]$ . Are the variables  $P_{i,j}$  and  $P_{k,l}$  still independent? What is the joint distribution  $\mathbf{P}(P_{1,1}, \dots, P_{4,4})$  now? Redo the calculation for the probabilities of pits in  $[1,3]$  and  $[2,2]$ .

**13.24** Redo the probability calculation for pits in  $[1,3]$  and  $[2,2]$ , assuming that each square contains a pit with probability 0.01, independent of the other squares. What can you say about the relative performance of a logical versus a probabilistic agent in this case?

**13.25** Implement a hybrid probabilistic agent for the wumpus world, based on the hybrid agent in Figure 7.20 and the probabilistic inference procedure outlined in this chapter.



# 14 PROBABILISTIC REASONING

*In which we explain how to build network models to reason under uncertainty according to the laws of probability theory.*

Chapter 13 introduced the basic elements of probability theory and noted the importance of independence and conditional independence relationships in simplifying probabilistic representations of the world. This chapter introduces a systematic way to represent such relationships explicitly in the form of **Bayesian networks**. We define the syntax and semantics of these networks and show how they can be used to capture uncertain knowledge in a natural and efficient way. We then show how probabilistic inference, although computationally intractable in the worst case, can be done efficiently in many practical situations. We also describe a variety of approximate inference algorithms that are often applicable when exact inference is infeasible. We explore ways in which probability theory can be applied to worlds with objects and relations—that is, to *first-order*, as opposed to *propositional*, representations. Finally, we survey alternative approaches to uncertain reasoning.

## 14.1 REPRESENTING KNOWLEDGE IN AN UNCERTAIN DOMAIN

In Chapter 13, we saw that the full joint probability distribution can answer any question about the domain, but can become intractably large as the number of variables grows. Furthermore, specifying probabilities for possible worlds one by one is unnatural and tedious.

We also saw that independence and conditional independence relationships among variables can greatly reduce the number of probabilities that need to be specified in order to define the full joint distribution. This section introduces a data structure called a **Bayesian network**<sup>1</sup> to represent the dependencies among variables. Bayesian networks can represent essentially *any* full joint probability distribution and in many cases can do so very concisely.

BAYESIAN NETWORK

<sup>1</sup> This is the most common name, but there are many synonyms, including **belief network**, **probabilistic network**, **causal network**, and **knowledge map**. In statistics, the term **graphical model** refers to a somewhat broader class that includes Bayesian networks. An extension of Bayesian networks called a **decision network** or **influence diagram** is covered in Chapter 16.

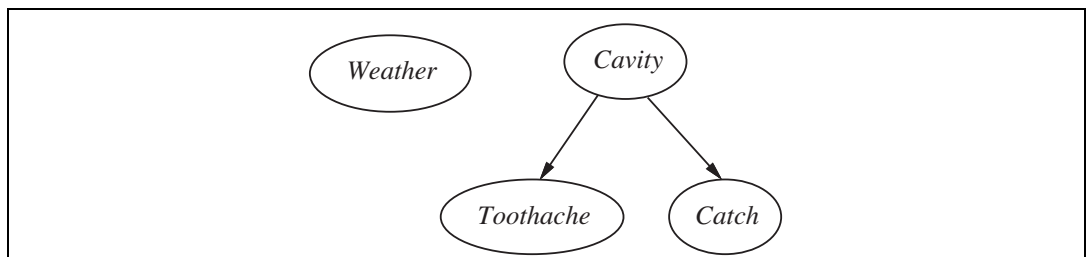
A Bayesian network is a directed graph in which each node is annotated with quantitative probability information. The full specification is as follows:

1. Each node corresponds to a random variable, which may be discrete or continuous.
2. A set of directed links or arrows connects pairs of nodes. If there is an arrow from node  $X$  to node  $Y$ ,  $X$  is said to be a *parent* of  $Y$ . The graph has no directed cycles (and hence is a directed acyclic graph, or DAG).
3. Each node  $X_i$  has a conditional probability distribution  $\mathbf{P}(X_i \mid \text{Parents}(X_i))$  that quantifies the effect of the parents on the node.

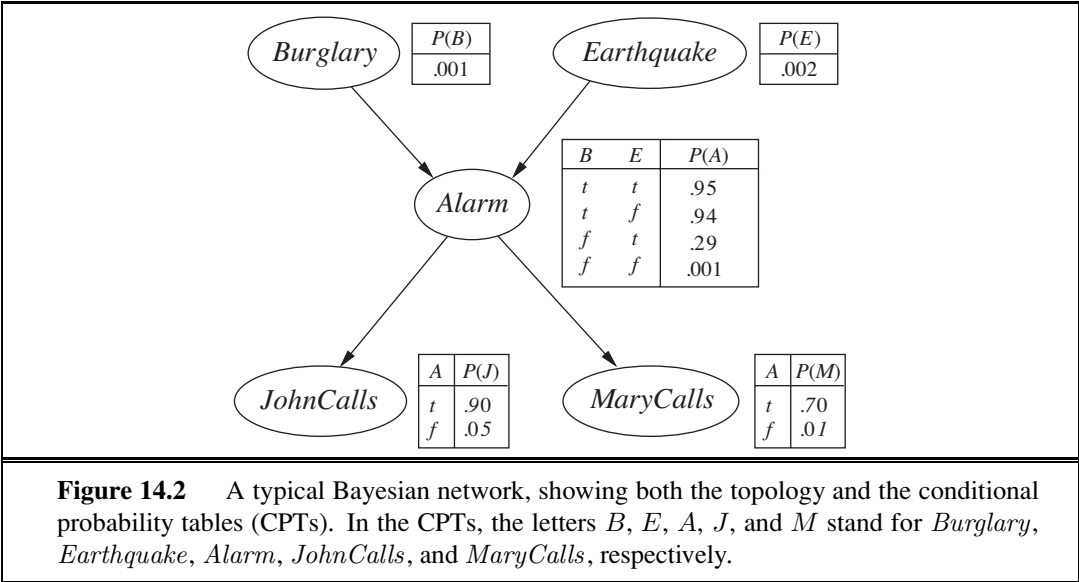
The topology of the network—the set of nodes and links—specifies the conditional independence relationships that hold in the domain, in a way that will be made precise shortly. The *intuitive* meaning of an arrow is typically that  $X$  has a *direct influence* on  $Y$ , which suggests that causes should be parents of effects. It is usually easy for a domain expert to decide what direct influences exist in the domain—much easier, in fact, than actually specifying the probabilities themselves. Once the topology of the Bayesian network is laid out, we need only specify a conditional probability distribution for each variable, given its parents. We will see that the combination of the topology and the conditional distributions suffices to specify (implicitly) the full joint distribution for all the variables.

Recall the simple world described in Chapter 13, consisting of the variables *Toothache*, *Cavity*, *Catch*, and *Weather*. We argued that *Weather* is independent of the other variables; furthermore, we argued that *Toothache* and *Catch* are conditionally independent, given *Cavity*. These relationships are represented by the Bayesian network structure shown in Figure 14.1. Formally, the conditional independence of *Toothache* and *Catch*, given *Cavity*, is indicated by the *absence* of a link between *Toothache* and *Catch*. Intuitively, the network represents the fact that *Cavity* is a direct cause of *Toothache* and *Catch*, whereas no direct causal relationship exists between *Toothache* and *Catch*.

Now consider the following example, which is just a little more complex. You have a new burglar alarm installed at home. It is fairly reliable at detecting a burglary, but also responds on occasion to minor earthquakes. (This example is due to Judea Pearl, a resident of Los Angeles—hence the acute interest in earthquakes.) You also have two neighbors, John and Mary, who have promised to call you at work when they hear the alarm. John nearly always calls when he hears the alarm, but sometimes confuses the telephone ringing with



**Figure 14.1** A simple Bayesian network in which *Weather* is independent of the other three variables and *Toothache* and *Catch* are conditionally independent, given *Cavity*.



the alarm and calls then, too. Mary, on the other hand, likes rather loud music and often misses the alarm altogether. Given the evidence of who has or has not called, we would like to estimate the probability of a burglary.

A Bayesian network for this domain appears in Figure 14.2. The network structure shows that burglary and earthquakes directly affect the probability of the alarm’s going off, but whether John and Mary call depends only on the alarm. The network thus represents our assumptions that they do not perceive burglaries directly, they do not notice minor earthquakes, and they do not confer before calling.

The conditional distributions in Figure 14.2 are shown as a **conditional probability table**, or CPT. (This form of table can be used for discrete variables; other representations, including those suitable for continuous variables, are described in Section 14.2.) Each row in a CPT contains the conditional probability of each node value for a **conditioning case**. A conditioning case is just a possible combination of values for the parent nodes—a miniature possible world, if you like. Each row must sum to 1, because the entries represent an exhaustive set of cases for the variable. For Boolean variables, once you know that the probability of a true value is  $p$ , the probability of false must be  $1 - p$ , so we often omit the second number, as in Figure 14.2. In general, a table for a Boolean variable with  $k$  Boolean parents contains  $2^k$  independently specifiable probabilities. A node with no parents has only one row, representing the prior probabilities of each possible value of the variable.

Notice that the network does not have nodes corresponding to Mary’s currently listening to loud music or to the telephone ringing and confusing John. These factors are summarized in the uncertainty associated with the links from *Alarm* to *JohnCalls* and *MaryCalls*. This shows both laziness and ignorance in operation: it would be a lot of work to find out why those factors would be more or less likely in any particular case, and we have no reasonable way to obtain the relevant information anyway. The probabilities actually summarize a *potentially*

CONDITIONAL  
PROBABILITY TABLE

CONDITIONING CASE

*infinite* set of circumstances in which the alarm might fail to go off (high humidity, power failure, dead battery, cut wires, a dead mouse stuck inside the bell, etc.) or John or Mary might fail to call and report it (out to lunch, on vacation, temporarily deaf, passing helicopter, etc.). In this way, a small agent can cope with a very large world, at least approximately. The degree of approximation can be improved if we introduce additional relevant information.

## 14.2 THE SEMANTICS OF BAYESIAN NETWORKS

The previous section described what a network is, but not what it means. There are two ways in which one can understand the semantics of Bayesian networks. The first is to see the network as a representation of the joint probability distribution. The second is to view it as an encoding of a collection of conditional independence statements. The two views are equivalent, but the first turns out to be helpful in understanding how to *construct* networks, whereas the second is helpful in designing inference procedures.

### 14.2.1 Representing the full joint distribution

Viewed as a piece of “syntax,” a Bayesian network is a directed acyclic graph with some numeric parameters attached to each node. One way to define what the network means—its semantics—is to define the way in which it represents a specific joint distribution over all the variables. To do this, we first need to retract (temporarily) what we said earlier about the parameters associated with each node. We said that those parameters correspond to conditional probabilities  $\mathbf{P}(X_i | \text{Parents}(X_i))$ ; this is a true statement, but until we assign semantics to the network as a whole, we should think of them just as numbers  $\theta(X_i | \text{Parents}(X_i))$ .

A generic entry in the joint distribution is the probability of a conjunction of particular assignments to each variable, such as  $P(X_1 = x_1 \wedge \dots \wedge X_n = x_n)$ . We use the notation  $P(x_1, \dots, x_n)$  as an abbreviation for this. The value of this entry is given by the formula

$$P(x_1, \dots, x_n) = \prod_{i=1}^n \theta(x_i | \text{parents}(X_i)), \quad (14.1)$$

where  $\text{parents}(X_i)$  denotes the values of  $\text{Parents}(X_i)$  that appear in  $x_1, \dots, x_n$ . Thus, each entry in the joint distribution is represented by the product of the appropriate elements of the conditional probability tables (CPTs) in the Bayesian network.

From this definition, it is easy to prove that the parameters  $\theta(X_i | \text{Parents}(X_i))$  are exactly the conditional probabilities  $\mathbf{P}(X_i | \text{Parents}(X_i))$  implied by the joint distribution (see Exercise 14.2). Hence, we can rewrite Equation (14.1) as

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i)). \quad (14.2)$$

In other words, the tables we have been calling conditional probability tables really *are* conditional probability tables according to the semantics defined in Equation (14.1).

To illustrate this, we can calculate the probability that the alarm has sounded, but neither a burglary nor an earthquake has occurred, and both John and Mary call. We multiply entries

from the joint distribution (using single-letter names for the variables):

$$\begin{aligned} P(j, m, a, \neg b, \neg e) &= P(j | a)P(m | a)P(a | \neg b \wedge \neg e)P(\neg b)P(\neg e) \\ &= 0.90 \times 0.70 \times 0.001 \times 0.999 \times 0.998 = 0.000628 . \end{aligned}$$

Section 13.3 explained that the full joint distribution can be used to answer any query about the domain. If a Bayesian network is a representation of the joint distribution, then it too can be used to answer any query, by summing all the relevant joint entries. Section 14.4 explains how to do this, but also describes methods that are much more efficient.

### A method for constructing Bayesian networks

Equation (14.2) defines what a given Bayesian network means. The next step is to explain how to *construct* a Bayesian network in such a way that the resulting joint distribution is a good representation of a given domain. We will now show that Equation (14.2) implies certain conditional independence relationships that can be used to guide the knowledge engineer in constructing the topology of the network. First, we rewrite the entries in the joint distribution in terms of conditional probability, using the product rule (see page 486):

$$P(x_1, \dots, x_n) = P(x_n | x_{n-1}, \dots, x_1)P(x_{n-1}, \dots, x_1) .$$

Then we repeat the process, reducing each conjunctive probability to a conditional probability and a smaller conjunction. We end up with one big product:

$$\begin{aligned} P(x_1, \dots, x_n) &= P(x_n | x_{n-1}, \dots, x_1)P(x_{n-1} | x_{n-2}, \dots, x_1) \cdots P(x_2 | x_1)P(x_1) \\ &= \prod_{i=1}^n P(x_i | x_{i-1}, \dots, x_1) . \end{aligned}$$

CHAIN RULE

This identity is called the **chain rule**. It holds for any set of random variables. Comparing it with Equation (14.2), we see that the specification of the joint distribution is equivalent to the general assertion that, for every variable  $X_i$  in the network,

$$\mathbf{P}(X_i | X_{i-1}, \dots, X_1) = \mathbf{P}(X_i | \text{Parents}(X_i)) , \quad (14.3)$$

provided that  $\text{Parents}(X_i) \subseteq \{X_{i-1}, \dots, X_1\}$ . This last condition is satisfied by numbering the nodes in a way that is consistent with the partial order implicit in the graph structure.

What Equation (14.3) says is that the Bayesian network is a correct representation of the domain only if each node is conditionally independent of its other predecessors in the node ordering, given its parents. We can satisfy this condition with this methodology:

1. *Nodes*: First determine the set of variables that are required to model the domain. Now order them,  $\{X_1, \dots, X_n\}$ . Any order will work, but the resulting network will be more compact if the variables are ordered such that causes precede effects.
2. *Links*: For  $i = 1$  to  $n$  do:
  - Choose, from  $X_1, \dots, X_{i-1}$ , a minimal set of parents for  $X_i$ , such that Equation (14.3) is satisfied.
  - For each parent insert a link from the parent to  $X_i$ .
  - CPTs: Write down the conditional probability table,  $\mathbf{P}(X_i | \text{Parents}(X_i))$ .



Intuitively, the parents of node  $X_i$  should contain all those nodes in  $X_1, \dots, X_{i-1}$  that *directly influence*  $X_i$ . For example, suppose we have completed the network in Figure 14.2 except for the choice of parents for *MaryCalls*. *MaryCalls* is certainly influenced by whether there is a *Burglary* or an *Earthquake*, but not *directly* influenced. Intuitively, our knowledge of the domain tells us that these events influence Mary's calling behavior only through their effect on the alarm. Also, given the state of the alarm, whether John calls has no influence on Mary's calling. Formally speaking, we believe that the following conditional independence statement holds:

$$\mathbf{P}(\text{MaryCalls} \mid \text{JohnCalls}, \text{Alarm}, \text{Earthquake}, \text{Burglary}) = \mathbf{P}(\text{MaryCalls} \mid \text{Alarm}) .$$

Thus, *Alarm* will be the only parent node for *MaryCalls*.

Because each node is connected only to earlier nodes, this construction method guarantees that the network is acyclic. Another important property of Bayesian networks is that they contain no redundant probability values. If there is no redundancy, then there is no chance for inconsistency: *it is impossible for the knowledge engineer or domain expert to create a Bayesian network that violates the axioms of probability.*

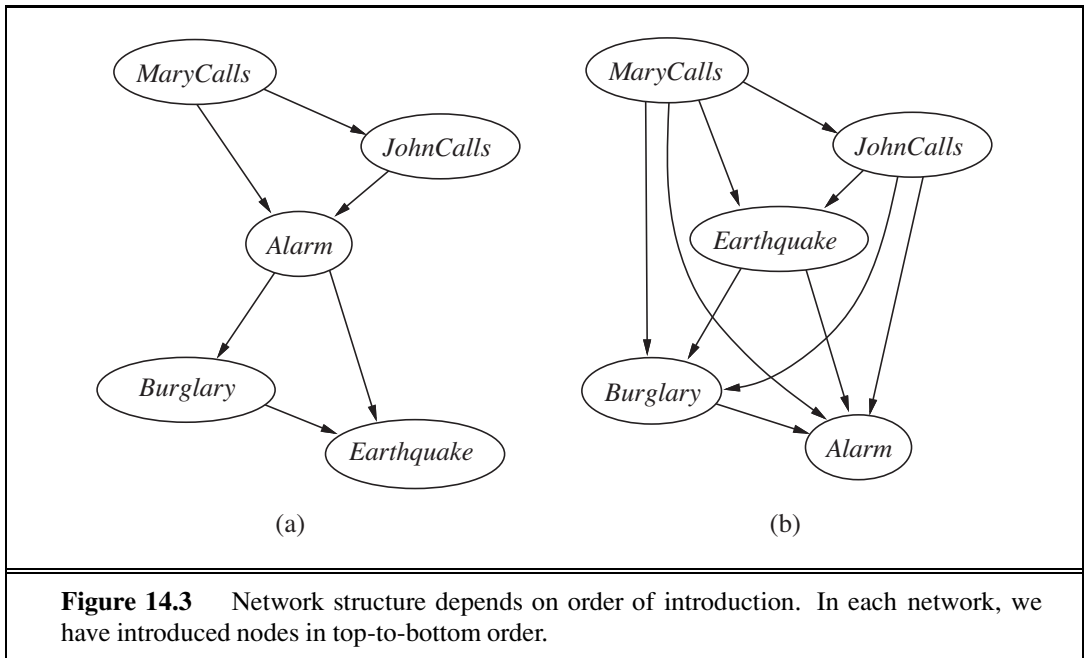


### Compactness and node ordering

As well as being a complete and nonredundant representation of the domain, a Bayesian network can often be far more *compact* than the full joint distribution. This property is what makes it feasible to handle domains with many variables. The compactness of Bayesian networks is an example of a general property of **locally structured** (also called **sparse**) systems. In a locally structured system, each subcomponent interacts directly with only a bounded number of other components, regardless of the total number of components. Local structure is usually associated with linear rather than exponential growth in complexity. In the case of Bayesian networks, it is reasonable to suppose that in most domains each random variable is directly influenced by at most  $k$  others, for some constant  $k$ . If we assume  $n$  Boolean variables for simplicity, then the amount of information needed to specify each conditional probability table will be at most  $2^k$  numbers, and the complete network can be specified by  $n2^k$  numbers. In contrast, the joint distribution contains  $2^n$  numbers. To make this concrete, suppose we have  $n = 30$  nodes, each with five parents ( $k = 5$ ). Then the Bayesian network requires 960 numbers, but the full joint distribution requires over a billion.

There are domains in which each variable can be influenced directly by all the others, so that the network is fully connected. Then specifying the conditional probability tables requires the same amount of information as specifying the joint distribution. In some domains, there will be slight dependencies that should strictly be included by adding a new link. But if these dependencies are tenuous, then it may not be worth the additional complexity in the network for the small gain in accuracy. For example, one might object to our burglary network on the grounds that if there is an earthquake, then John and Mary would not call even if they heard the alarm, because they assume that the earthquake is the cause. Whether to add the link from *Earthquake* to *JohnCalls* and *MaryCalls* (and thus enlarge the tables) depends on comparing the importance of getting more accurate probabilities with the cost of specifying the extra information.

LOCALLY  
STRUCTURED  
SPARSE



Even in a locally structured domain, we will get a compact Bayesian network only if we choose the node ordering well. What happens if we happen to choose the wrong order? Consider the burglary example again. Suppose we decide to add the nodes in the order *MaryCalls*, *JohnCalls*, *Alarm*, *Burglary*, *Earthquake*. We then get the somewhat more complicated network shown in Figure 14.3(a). The process goes as follows:

- Adding *MaryCalls*: No parents.
- Adding *JohnCalls*: If Mary calls, that probably means the alarm has gone off, which of course would make it more likely that John calls. Therefore, *JohnCalls* needs *MaryCalls* as a parent.
- Adding *Alarm*: Clearly, if both call, it is more likely that the alarm has gone off than if just one or neither calls, so we need both *MaryCalls* and *JohnCalls* as parents.
- Adding *Burglary*: If we know the alarm state, then the call from John or Mary might give us information about our phone ringing or Mary's music, but not about burglary:

$$\mathbf{P}(\text{Burglary} \mid \text{Alarm}, \text{JohnCalls}, \text{MaryCalls}) = \mathbf{P}(\text{Burglary} \mid \text{Alarm}) .$$

Hence we need just *Alarm* as parent.

- Adding *Earthquake*: If the alarm is on, it is more likely that there has been an earthquake. (The alarm is an earthquake detector of sorts.) But if we know that there has been a burglary, then that explains the alarm, and the probability of an earthquake would be only slightly above normal. Hence, we need both *Alarm* and *Burglary* as parents.

The resulting network has two more links than the original network in Figure 14.2 and requires three more probabilities to be specified. What's worse, some of the links represent tenuous relationships that require difficult and unnatural probability judgments, such as as-





sessing the probability of *Earthquake*, given *Burglary* and *Alarm*. This phenomenon is quite general and is related to the distinction between **causal** and **diagnostic** models introduced in Section 13.5.1 (see also Exercise 8.13). If we try to build a diagnostic model with links from symptoms to causes (as from *MaryCalls* to *Alarm* or *Alarm* to *Burglary*), we end up having to specify additional dependencies between otherwise independent causes (and often between separately occurring symptoms as well). *If we stick to a causal model, we end up having to specify fewer numbers, and the numbers will often be easier to come up with.* In the domain of medicine, for example, it has been shown by Tversky and Kahneman (1982) that expert physicians prefer to give probability judgments for causal rules rather than for diagnostic ones.

Figure 14.3(b) shows a very bad node ordering: *MaryCalls*, *JohnCalls*, *Earthquake*, *Burglary*, *Alarm*. This network requires 31 distinct probabilities to be specified—exactly the same number as the full joint distribution. It is important to realize, however, that any of the three networks can represent *exactly the same joint distribution*. The last two versions simply fail to represent all the conditional independence relationships and hence end up specifying a lot of unnecessary numbers instead.

## 14.2.2 Conditional independence relations in Bayesian networks

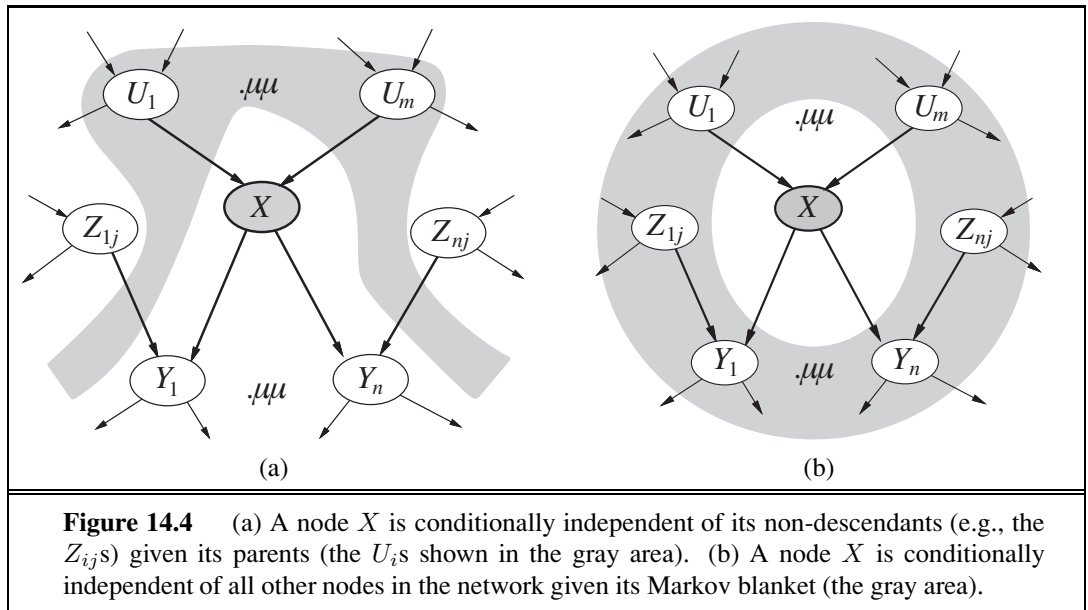
We have provided a “numerical” semantics for Bayesian networks in terms of the representation of the full joint distribution, as in Equation (14.2). Using this semantics to derive a method for constructing Bayesian networks, we were led to the consequence that a node is conditionally independent of its other predecessors, given its parents. It turns out that we can also go in the other direction. We can start from a “topological” semantics that specifies the conditional independence relationships encoded by the graph structure, and from this we can derive the “numerical” semantics. The topological semantics<sup>2</sup> specifies that each variable is conditionally independent of its non-**descendants**, given its parents. For example, in Figure 14.2, *JohnCalls* is independent of *Burglary*, *Earthquake*, and *MaryCalls* given the value of *Alarm*. The definition is illustrated in Figure 14.4(a). From these conditional independence assertions and the interpretation of the network parameters  $\theta(X_i | \text{Parents}(X_i))$  as specifications of conditional probabilities  $\mathbf{P}(X_i | \text{Parents}(X_i))$ , the full joint distribution given in Equation (14.2) can be reconstructed. In this sense, the “numerical” semantics and the “topological” semantics are equivalent.

Another important independence property is implied by the topological semantics: a node is conditionally independent of all other nodes in the network, given its parents, children, and children’s parents—that is, given its **Markov blanket**. (Exercise 14.7 asks you to prove this.) For example, *Burglary* is independent of *JohnCalls* and *MaryCalls*, given *Alarm* and *Earthquake*. This property is illustrated in Figure 14.4(b).

<sup>2</sup> There is also a general topological criterion called **d-separation** for deciding whether a set of nodes **X** is conditionally independent of another set **Y**, given a third set **Z**. The criterion is rather complicated and is not needed for deriving the algorithms in this chapter, so we omit it. Details may be found in Pearl (1988) or Darwiche (2009). Shachter (1998) gives a more intuitive method of ascertaining d-separation.

DESCENDANT

MARKOV BLANKET



### 14.3 EFFICIENT REPRESENTATION OF CONDITIONAL DISTRIBUTIONS

Even if the maximum number of parents  $k$  is smallish, filling in the CPT for a node requires up to  $O(2^k)$  numbers and perhaps a great deal of experience with all the possible conditioning cases. In fact, this is a worst-case scenario in which the relationship between the parents and the child is completely arbitrary. Usually, such relationships are describable by a **canonical distribution** that fits some standard pattern. In such cases, the complete table can be specified by naming the pattern and perhaps supplying a few parameters—much easier than supplying an exponential number of parameters.

The simplest example is provided by **deterministic nodes**. A deterministic node has its value specified exactly by the values of its parents, with no uncertainty. The relationship can be a logical one: for example, the relationship between the parent nodes *Canadian*, *US*, *Mexican* and the child node *NorthAmerican* is simply that the child is the disjunction of the parents. The relationship can also be numerical: for example, if the parent nodes are the prices of a particular model of car at several dealers and the child node is the price that a bargain hunter ends up paying, then the child node is the minimum of the parent values; or if the parent nodes are a lake's inflows (rivers, runoff, precipitation) and outflows (rivers, evaporation, seepage) and the child is the change in the water level of the lake, then the value of the child is the sum of the inflow parents minus the sum of the outflow parents.

Uncertain relationships can often be characterized by so-called **noisy** logical relationships. The standard example is the **noisy-OR** relation, which is a generalization of the logical OR. In propositional logic, we might say that *Fever* is true if and only if *Cold*, *Flu*, or *Malaria* is true. The noisy-OR model allows for uncertainty about the ability of each parent to cause the child to be true—the causal relationship between parent and child may be

CANONICAL  
DISTRIBUTION

DETERMINISTIC  
NODES

NOISY-OR

LEAK NODE

*inhibited*, and so a patient could have a cold, but not exhibit a fever. The model makes two assumptions. First, it assumes that all the possible causes are listed. (If some are missing, we can always add a so-called **leak node** that covers “miscellaneous causes.”) Second, it assumes that inhibition of each parent is independent of inhibition of any other parents: for example, whatever inhibits *Malaria* from causing a fever is independent of whatever inhibits *Flu* from causing a fever. Given these assumptions, *Fever* is *false* if and only if all its *true* parents are inhibited, and the probability of this is the product of the inhibition probabilities  $q$  for each parent. Let us suppose these individual inhibition probabilities are as follows:

$$\begin{aligned} q_{\text{cold}} &= P(\neg \text{fever} \mid \text{cold}, \neg \text{flu}, \neg \text{malaria}) = 0.6, \\ q_{\text{flu}} &= P(\neg \text{fever} \mid \neg \text{cold}, \text{flu}, \neg \text{malaria}) = 0.2, \\ q_{\text{malaria}} &= P(\neg \text{fever} \mid \neg \text{cold}, \neg \text{flu}, \text{malaria}) = 0.1. \end{aligned}$$

Then, from this information and the noisy-OR assumptions, the entire CPT can be built. The general rule is that

$$P(x_i \mid \text{parents}(X_i)) = 1 - \prod_{\{j: X_j = \text{true}\}} q_j,$$

where the product is taken over the parents that are set to true for that row of the CPT. The following table illustrates this calculation:

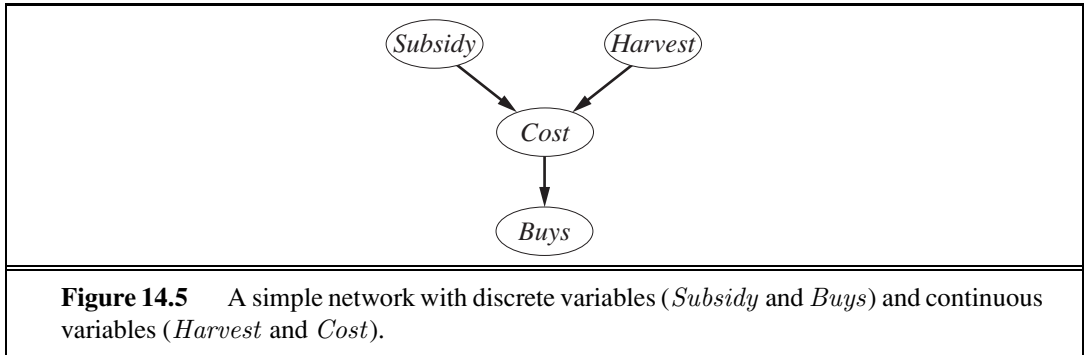
<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	$P(\text{Fever})$	$P(\neg \text{Fever})$
F	F	F	0.0	1.0
F	F	T	0.9	<b>0.1</b>
F	T	F	0.8	<b>0.2</b>
F	T	T	0.98	$0.02 = 0.2 \times 0.1$
T	F	F	0.4	<b>0.6</b>
T	F	T	0.94	$0.06 = 0.6 \times 0.1$
T	T	F	0.88	$0.12 = 0.6 \times 0.2$
T	T	T	0.988	$0.012 = 0.6 \times 0.2 \times 0.1$

In general, noisy logical relationships in which a variable depends on  $k$  parents can be described using  $O(k)$  parameters instead of  $O(2^k)$  for the full conditional probability table. This makes assessment and learning much easier. For example, the CPCS network (Pradhan *et al.*, 1994) uses noisy-OR and noisy-MAX distributions to model relationships among diseases and symptoms in internal medicine. With 448 nodes and 906 links, it requires only 8,254 values instead of 133,931,430 for a network with full CPTs.

### Bayesian nets with continuous variables

Many real-world problems involve continuous quantities, such as height, mass, temperature, and money; in fact, much of statistics deals with random variables whose domains are continuous. By definition, continuous variables have an infinite number of possible values, so it is impossible to specify conditional probabilities explicitly for each value. One possible way to handle continuous variables is to avoid them by using **discretization**—that is, dividing up the

DISCRETIZATION



possible values into a fixed set of intervals. For example, temperatures could be divided into ( $<0^\circ\text{C}$ ), ( $0^\circ\text{C}–100^\circ\text{C}$ ), and ( $>100^\circ\text{C}$ ). Discretization is sometimes an adequate solution, but often results in a considerable loss of accuracy and very large CPTs. The most common solution is to define standard families of probability density functions (see Appendix A) that are specified by a finite number of **parameters**. For example, a Gaussian (or normal) distribution  $N(\mu, \sigma^2)(x)$  has the mean  $\mu$  and the variance  $\sigma^2$  as parameters. Yet another solution—sometimes called a **nonparametric** representation—is to define the conditional distribution implicitly with a collection of instances, each containing specific values of the parent and child variables. We explore this approach further in Chapter 18.

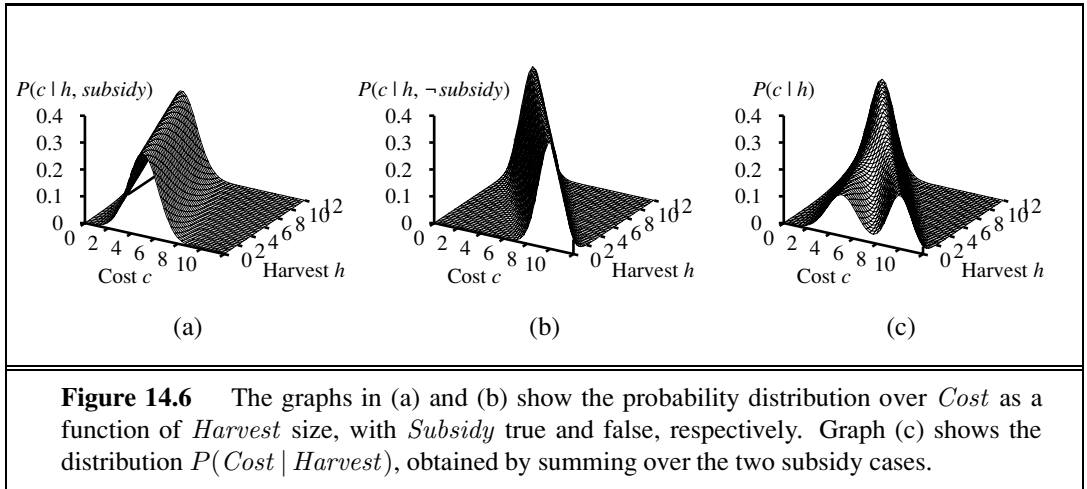
A network with both discrete and continuous variables is called a **hybrid Bayesian network**. To specify a hybrid network, we have to specify two new kinds of distributions: the conditional distribution for a continuous variable given discrete or continuous parents; and the conditional distribution for a discrete variable given continuous parents. Consider the simple example in Figure 14.5, in which a customer buys some fruit depending on its cost, which depends in turn on the size of the harvest and whether the government’s subsidy scheme is operating. The variable *Cost* is continuous and has continuous and discrete parents; the variable *Buys* is discrete and has a continuous parent.

For the *Cost* variable, we need to specify  $\mathbf{P}(\text{Cost} \mid \text{Harvest}, \text{Subsidy})$ . The discrete parent is handled by enumeration—that is, by specifying both  $P(\text{Cost} \mid \text{Harvest}, \text{subsidy})$  and  $P(\text{Cost} \mid \text{Harvest}, \neg\text{subsidy})$ . To handle *Harvest*, we specify how the distribution over the cost  $c$  depends on the continuous value  $h$  of *Harvest*. In other words, we specify the *parameters* of the cost distribution as a function of  $h$ . The most common choice is the **linear Gaussian** distribution, in which the child has a Gaussian distribution whose mean  $\mu$  varies linearly with the value of the parent and whose standard deviation  $\sigma$  is fixed. We need two distributions, one for *subsidy* and one for  $\neg\text{subsidy}$ , with different parameters:

$$P(c \mid h, \text{subsidy}) = N(a_t h + b_t, \sigma_t^2)(c) = \frac{1}{\sigma_t \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{c - (a_t h + b_t)}{\sigma_t} \right)^2}$$

$$P(c \mid h, \neg\text{subsidy}) = N(a_f h + b_f, \sigma_f^2)(c) = \frac{1}{\sigma_f \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{c - (a_f h + b_f)}{\sigma_f} \right)^2}.$$

For this example, then, the conditional distribution for *Cost* is specified by naming the linear Gaussian distribution and providing the parameters  $a_t, b_t, \sigma_t, a_f, b_f$ , and  $\sigma_f$ . Figures 14.6(a)



and (b) show these two relationships. Notice that in each case the slope is negative, because cost decreases as supply increases. (Of course, the assumption of linearity implies that the cost becomes negative at some point; the linear model is reasonable only if the harvest size is limited to a narrow range.) Figure 14.6(c) shows the distribution  $P(c | h)$ , averaging over the two possible values of *Subsidy* and assuming that each has prior probability 0.5. This shows that even with very simple models, quite interesting distributions can be represented.

The linear Gaussian conditional distribution has some special properties. A network containing only continuous variables with linear Gaussian distributions has a joint distribution that is a multivariate Gaussian distribution (see Appendix A) over all the variables (Exercise 14.9). Furthermore, the posterior distribution given any evidence also has this property.<sup>3</sup> When discrete variables are added as parents (not as children) of continuous variables, the network defines a **conditional Gaussian**, or CG, distribution: given any assignment to the discrete variables, the distribution over the continuous variables is a multivariate Gaussian.

Now we turn to the distributions for discrete variables with continuous parents. Consider, for example, the *Buys* node in Figure 14.5. It seems reasonable to assume that the customer will buy if the cost is low and will not buy if it is high and that the probability of buying varies smoothly in some intermediate region. In other words, the conditional distribution is like a “soft” threshold function. One way to make soft thresholds is to use the *integral* of the standard normal distribution:

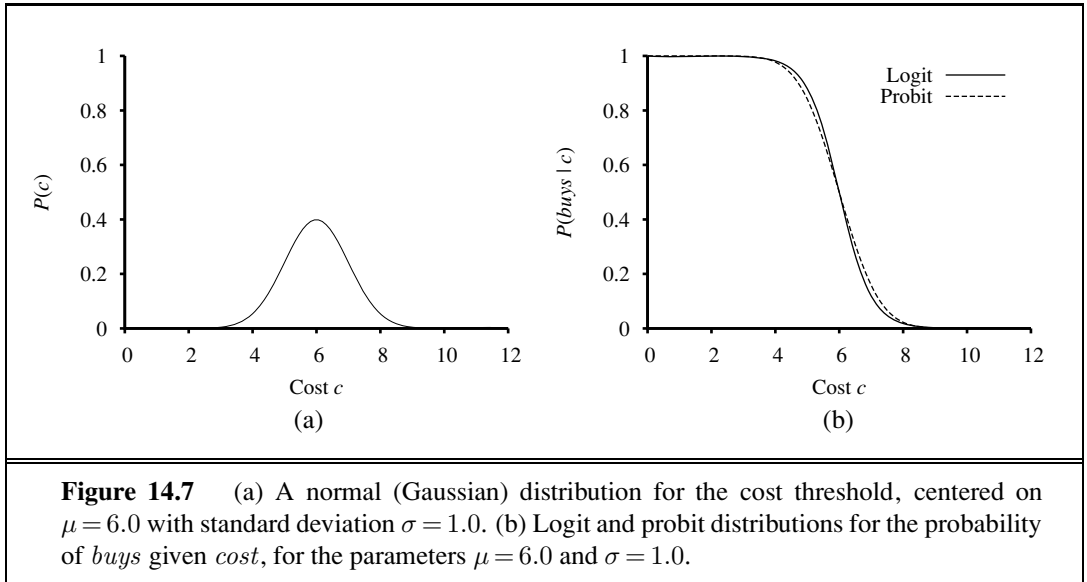
$$\Phi(x) = \int_{-\infty}^x N(0, 1)(x) dx .$$

Then the probability of *Buys* given *Cost* might be

$$P(\text{buys} | \text{Cost} = c) = \Phi((-c + \mu)/\sigma) ,$$

which means that the cost threshold occurs around  $\mu$ , the width of the threshold region is proportional to  $\sigma$ , and the probability of buying decreases as cost increases. This **probit distribution**

<sup>3</sup> It follows that inference in linear Gaussian networks takes only  $O(n^3)$  time in the worst case, regardless of the network topology. In Section 14.4, we see that inference for networks of discrete variables is NP-hard.



PROBIT  
DISTRIBUTION

**bution** (pronounced “pro-bit” and short for “probability unit”) is illustrated in Figure 14.7(a). The form can be justified by proposing that the underlying decision process has a hard threshold, but that the precise location of the threshold is subject to random Gaussian noise.

LOGIT DISTRIBUTION  
LOGISTIC FUNCTION

An alternative to the probit model is the **logit distribution** (pronounced “low-jit”). It uses the **logistic function**  $1/(1 + e^{-x})$  to produce a soft threshold:

$$P(buys | Cost = c) = \frac{1}{1 + \exp(-2\frac{c - \mu}{\sigma})}.$$

This is illustrated in Figure 14.7(b). The two distributions look similar, but the logit actually has much longer “tails.” The probit is often a better fit to real situations, but the logit is sometimes easier to deal with mathematically. It is used widely in neural networks (Chapter 20). Both probit and logit can be generalized to handle multiple continuous parents by taking a linear combination of the parent values.

## 14.4 EXACT INFERENCE IN BAYESIAN NETWORKS

EVENT

The basic task for any probabilistic inference system is to compute the posterior probability distribution for a set of **query variables**, given some observed **event**—that is, some assignment of values to a set of **evidence variables**. To simplify the presentation, we will consider only one query variable at a time; the algorithms can easily be extended to queries with multiple variables. We will use the notation from Chapter 13:  $X$  denotes the query variable;  $\mathbf{E}$  denotes the set of evidence variables  $E_1, \dots, E_m$ , and  $\mathbf{e}$  is a particular observed event;  $\mathbf{Y}$  will denote the nonevidence, nonquery variables  $Y_1, \dots, Y_l$  (called the **hidden variables**). Thus, the complete set of variables is  $\mathbf{X} = \{X\} \cup \mathbf{E} \cup \mathbf{Y}$ . A typical query asks for the posterior probability distribution  $\mathbf{P}(X | \mathbf{e})$ .

HIDDEN VARIABLE

In the burglary network, we might observe the event in which  $JohnCalls = true$  and  $MaryCalls = true$ . We could then ask for, say, the probability that a burglary has occurred:

$$\mathbf{P}(\text{Burglary} \mid JohnCalls = true, MaryCalls = true) = \langle 0.284, 0.716 \rangle .$$

In this section we discuss exact algorithms for computing posterior probabilities and will consider the complexity of this task. It turns out that the general case is intractable, so Section 14.5 covers methods for approximate inference.

### 14.4.1 Inference by enumeration

Chapter 13 explained that any conditional probability can be computed by summing terms from the full joint distribution. More specifically, a query  $\mathbf{P}(X \mid \mathbf{e})$  can be answered using Equation (13.9), which we repeat here for convenience:

$$\mathbf{P}(X \mid \mathbf{e}) = \alpha \mathbf{P}(X, \mathbf{e}) = \alpha \sum_{\mathbf{y}} \mathbf{P}(X, \mathbf{e}, \mathbf{y}) .$$

Now, a Bayesian network gives a complete representation of the full joint distribution. More specifically, Equation (14.2) on page 513 shows that the terms  $P(x, \mathbf{e}, \mathbf{y})$  in the joint distribution can be written as products of conditional probabilities from the network. Therefore, *a query can be answered using a Bayesian network by computing sums of products of conditional probabilities from the network.*

Consider the query  $\mathbf{P}(\text{Burglary} \mid JohnCalls = true, MaryCalls = true)$ . The hidden variables for this query are *Earthquake* and *Alarm*. From Equation (13.9), using initial letters for the variables to shorten the expressions, we have<sup>4</sup>

$$\mathbf{P}(B \mid j, m) = \alpha \mathbf{P}(B, j, m) = \alpha \sum_e \sum_a \mathbf{P}(B, j, m, e, a) .$$

The semantics of Bayesian networks (Equation (14.2)) then gives us an expression in terms of CPT entries. For simplicity, we do this just for  $Burglary = true$ :

$$P(b \mid j, m) = \alpha \sum_e \sum_a P(b)P(e)P(a \mid b, e)P(j \mid a)P(m \mid a) .$$

To compute this expression, we have to add four terms, each computed by multiplying five numbers. In the worst case, where we have to sum out almost all the variables, the complexity of the algorithm for a network with  $n$  Boolean variables is  $O(n2^n)$ .

An improvement can be obtained from the following simple observations: the  $P(b)$  term is a constant and can be moved outside the summations over  $a$  and  $e$ , and the  $P(e)$  term can be moved outside the summation over  $a$ . Hence, we have

$$P(b \mid j, m) = \alpha P(b) \sum_e P(e) \sum_a P(a \mid b, e)P(j \mid a)P(m \mid a) . \quad (14.4)$$

This expression can be evaluated by looping through the variables in order, multiplying CPT entries as we go. For each summation, we also need to loop over the variable's possible

<sup>4</sup> An expression such as  $\sum_e P(a, e)$  means to sum  $P(A = a, E = e)$  for all possible values of  $e$ . When  $E$  is Boolean, there is an ambiguity in that  $P(e)$  is used to mean both  $P(E = true)$  and  $P(E = e)$ , but it should be clear from context which is intended; in particular, in the context of a sum the latter is intended.



values. The structure of this computation is shown in Figure 14.8. Using the numbers from Figure 14.2, we obtain  $P(b | j, m) = \alpha \times 0.00059224$ . The corresponding computation for  $\neg b$  yields  $\alpha \times 0.0014919$ ; hence,

$$\mathbf{P}(B | j, m) = \alpha \langle 0.00059224, 0.0014919 \rangle \approx \langle 0.284, 0.716 \rangle.$$

That is, the chance of a burglary, given calls from both neighbors, is about 28%.

The evaluation process for the expression in Equation (14.4) is shown as an expression tree in Figure 14.8. The ENUMERATION-ASK algorithm in Figure 14.9 evaluates such trees using depth-first recursion. The algorithm is very similar in structure to the backtracking algorithm for solving CSPs (Figure 6.5) and the DPLL algorithm for satisfiability (Figure 7.17).

The space complexity of ENUMERATION-ASK is only linear in the number of variables: the algorithm sums over the full joint distribution without ever constructing it explicitly. Unfortunately, its time complexity for a network with  $n$  Boolean variables is always  $O(2^n)$ —better than the  $O(n 2^n)$  for the simple approach described earlier, but still rather grim.

Note that the tree in Figure 14.8 makes explicit the *repeated subexpressions* evaluated by the algorithm. The products  $P(j | a)P(m | a)$  and  $P(j | \neg a)P(m | \neg a)$  are computed twice, once for each value of  $e$ . The next section describes a general method that avoids such wasted computations.

### 14.4.2 The variable elimination algorithm

The enumeration algorithm can be improved substantially by eliminating repeated calculations of the kind illustrated in Figure 14.8. The idea is simple: do the calculation once and save the results for later use. This is a form of dynamic programming. There are several versions of this approach; we present the **variable elimination** algorithm, which is the simplest. Variable elimination works by evaluating expressions such as Equation (14.4) in *right-to-left* order (that is, *bottom up* in Figure 14.8). Intermediate results are stored, and summations over each variable are done only for those portions of the expression that depend on the variable.

Let us illustrate this process for the burglary network. We evaluate the expression

$$\mathbf{P}(B | j, m) = \alpha \underbrace{\mathbf{P}(B)}_{\mathbf{f}_1(B)} \sum_e \underbrace{P(e)}_{\mathbf{f}_2(E)} \sum_a \underbrace{\mathbf{P}(a | B, e)}_{\mathbf{f}_3(A, B, E)} \underbrace{P(j | a)}_{\mathbf{f}_4(A)} \underbrace{P(m | a)}_{\mathbf{f}_5(A)}.$$

Notice that we have annotated each part of the expression with the name of the corresponding **factor**; each factor is a matrix indexed by the values of its argument variables. For example, the factors  $\mathbf{f}_4(A)$  and  $\mathbf{f}_5(A)$  corresponding to  $P(j | a)$  and  $P(m | a)$  depend just on  $A$  because  $J$  and  $M$  are fixed by the query. They are therefore two-element vectors:

$$\mathbf{f}_4(A) = \begin{pmatrix} P(j | a) \\ P(j | \neg a) \end{pmatrix} = \begin{pmatrix} 0.90 \\ 0.05 \end{pmatrix} \quad \mathbf{f}_5(A) = \begin{pmatrix} P(m | a) \\ P(m | \neg a) \end{pmatrix} = \begin{pmatrix} 0.70 \\ 0.01 \end{pmatrix}.$$

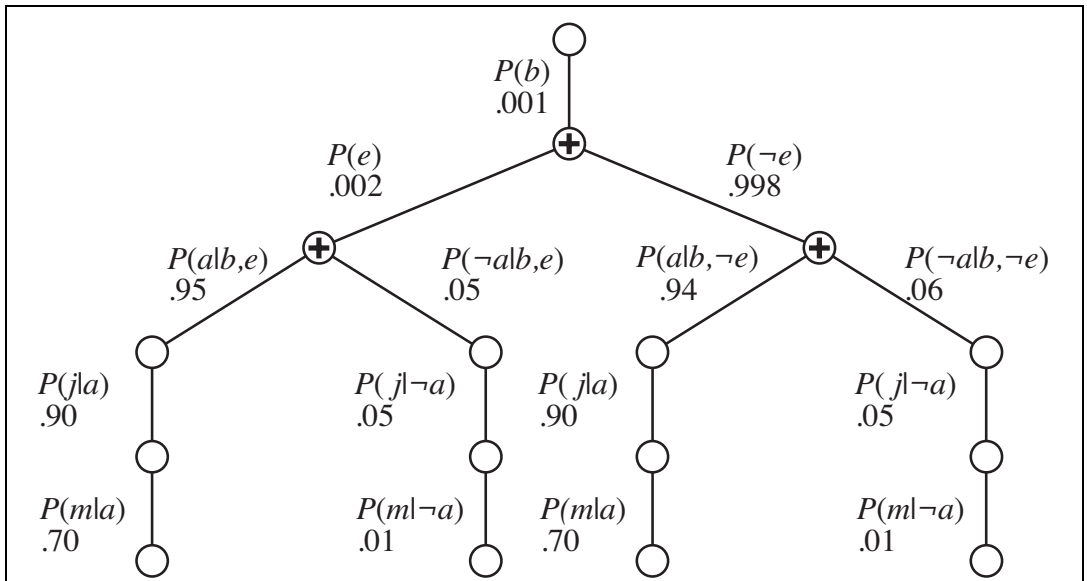
$\mathbf{f}_3(A, B, E)$  will be a  $2 \times 2 \times 2$  matrix, which is hard to show on the printed page. (The “first” element is given by  $P(a | b, e) = 0.95$  and the “last” by  $P(\neg a | \neg b, \neg e) = 0.999$ .) In terms of factors, the query expression is written as

$$\mathbf{P}(B | j, m) = \alpha \mathbf{f}_1(B) \times \sum_e \mathbf{f}_2(E) \times \sum_a \mathbf{f}_3(A, B, E) \times \mathbf{f}_4(A) \times \mathbf{f}_5(A)$$

VARIABLE  
ELIMINATION

FACTOR





**Figure 14.8** The structure of the expression shown in Equation (14.4). The evaluation proceeds top down, multiplying values along each path and summing at the “+” nodes. Notice the repetition of the paths for  $j$  and  $m$ .

```

function ENUMERATION-ASK( $X, \mathbf{e}, bn$ ) returns a distribution over  $X$ 
  inputs:  $X$ , the query variable
            $\mathbf{e}$ , observed values for variables  $\mathbf{E}$ 
            $bn$ , a Bayes net with variables  $\{X\} \cup \mathbf{E} \cup \mathbf{Y}$  /*  $\mathbf{Y} = \text{hidden variables}$  */

   $\mathbf{Q}(X) \leftarrow$  a distribution over  $X$ , initially empty
  for each value  $x_i$  of  $X$  do
     $\mathbf{Q}(x_i) \leftarrow$  ENUMERATE-ALL( $bn.VARS, \mathbf{e}_{x_i}$ )
    where  $\mathbf{e}_{x_i}$  is  $\mathbf{e}$  extended with  $X = x_i$ 
  return NORMALIZE( $\mathbf{Q}(X)$ )



---


function ENUMERATE-ALL( $vars, \mathbf{e}$ ) returns a real number
  if EMPTY?( $vars$ ) then return 1.0
   $Y \leftarrow$  FIRST( $vars$ )
  if  $Y$  has value  $y$  in  $\mathbf{e}$ 
    then return  $P(y \mid \text{parents}(Y)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}$ )
    else return  $\sum_y P(y \mid \text{parents}(Y)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}_y$ )
    where  $\mathbf{e}_y$  is  $\mathbf{e}$  extended with  $Y = y$ 

```

**Figure 14.9** The enumeration algorithm for answering queries on Bayesian networks.

where the “ $\times$ ” operator is not ordinary matrix multiplication but instead the **pointwise product** operation, to be described shortly.

The process of evaluation is a process of summing out variables (right to left) from pointwise products of factors to produce new factors, eventually yielding a factor that is the solution, i.e., the posterior distribution over the query variable. The steps are as follows:

- First, we sum out  $A$  from the product of  $\mathbf{f}_3$ ,  $\mathbf{f}_4$ , and  $\mathbf{f}_5$ . This gives us a new  $2 \times 2$  factor  $\mathbf{f}_6(B, E)$  whose indices range over just  $B$  and  $E$ :

$$\begin{aligned}\mathbf{f}_6(B, E) &= \sum_a \mathbf{f}_3(A, B, E) \times \mathbf{f}_4(A) \times \mathbf{f}_5(A) \\ &= (\mathbf{f}_3(a, B, E) \times \mathbf{f}_4(a) \times \mathbf{f}_5(a)) + (\mathbf{f}_3(\neg a, B, E) \times \mathbf{f}_4(\neg a) \times \mathbf{f}_5(\neg a)).\end{aligned}$$

Now we are left with the expression

$$\mathbf{P}(B \mid j, m) = \alpha \mathbf{f}_1(B) \times \sum_e \mathbf{f}_2(E) \times \mathbf{f}_6(B, E).$$

- Next, we sum out  $E$  from the product of  $\mathbf{f}_2$  and  $\mathbf{f}_6$ :

$$\begin{aligned}\mathbf{f}_7(B) &= \sum_e \mathbf{f}_2(E) \times \mathbf{f}_6(B, E) \\ &= \mathbf{f}_2(e) \times \mathbf{f}_6(B, e) + \mathbf{f}_2(\neg e) \times \mathbf{f}_6(B, \neg e).\end{aligned}$$

This leaves the expression

$$\mathbf{P}(B \mid j, m) = \alpha \mathbf{f}_1(B) \times \mathbf{f}_7(B)$$

which can be evaluated by taking the pointwise product and normalizing the result.

Examining this sequence, we see that two basic computational operations are required: pointwise product of a pair of factors, and summing out a variable from a product of factors. The next section describes each of these operations.

### Operations on factors

The pointwise product of two factors  $\mathbf{f}_1$  and  $\mathbf{f}_2$  yields a new factor  $\mathbf{f}$  whose variables are the *union* of the variables in  $\mathbf{f}_1$  and  $\mathbf{f}_2$  and whose elements are given by the product of the corresponding elements in the two factors. Suppose the two factors have variables  $Y_1, \dots, Y_k$  in common. Then we have

$$\mathbf{f}(X_1 \dots X_j, Y_1 \dots Y_k, Z_1 \dots Z_l) = \mathbf{f}_1(X_1 \dots X_j, Y_1 \dots Y_k) \mathbf{f}_2(Y_1 \dots Y_k, Z_1 \dots Z_l).$$

If all the variables are binary, then  $\mathbf{f}_1$  and  $\mathbf{f}_2$  have  $2^{j+k}$  and  $2^{k+l}$  entries, respectively, and the pointwise product has  $2^{j+k+l}$  entries. For example, given two factors  $\mathbf{f}_1(A, B)$  and  $\mathbf{f}_2(B, C)$ , the pointwise product  $\mathbf{f}_1 \times \mathbf{f}_2 = \mathbf{f}_3(A, B, C)$  has  $2^{1+1+1} = 8$  entries, as illustrated in Figure 14.10. Notice that the factor resulting from a pointwise product can contain more variables than any of the factors being multiplied and that the size of a factor is exponential in the number of variables. This is where both space and time complexity arise in the variable elimination algorithm.

$A$	$B$	$\mathbf{f}_1(A, B)$	$B$	$C$	$\mathbf{f}_2(B, C)$	$A$	$B$	$C$	$\mathbf{f}_3(A, B, C)$
T	T	.3	T	T	.2	T	T	T	$.3 \times .2 = .06$
T	F	.7	T	F	.8	T	T	F	$.3 \times .8 = .24$
F	T	.9	F	T	.6	T	F	T	$.7 \times .6 = .42$
F	F	.1	F	F	.4	T	F	F	$.7 \times .4 = .28$
						F	T	T	$.9 \times .2 = .18$
						F	T	F	$.9 \times .8 = .72$
						F	F	T	$.1 \times .6 = .06$
						F	F	F	$.1 \times .4 = .04$

**Figure 14.10** Illustrating pointwise multiplication:  $\mathbf{f}_1(A, B) \times \mathbf{f}_2(B, C) = \mathbf{f}_3(A, B, C)$ .

Summing out a variable from a product of factors is done by adding up the submatrices formed by fixing the variable to each of its values in turn. For example, to sum out  $A$  from  $\mathbf{f}_3(A, B, C)$ , we write

$$\begin{aligned} \mathbf{f}(B, C) &= \sum_a \mathbf{f}_3(A, B, C) = \mathbf{f}_3(a, B, C) + \mathbf{f}_3(\neg a, B, C) \\ &= \begin{pmatrix} .06 & .24 \\ .42 & .28 \end{pmatrix} + \begin{pmatrix} .18 & .72 \\ .06 & .04 \end{pmatrix} = \begin{pmatrix} .24 & .96 \\ .48 & .32 \end{pmatrix}. \end{aligned}$$

The only trick is to notice that any factor that does *not* depend on the variable to be summed out can be moved outside the summation. For example, if we were to sum out  $E$  first in the burglary network, the relevant part of the expression would be

$$\sum_e \mathbf{f}_2(E) \times \mathbf{f}_3(A, B, E) \times \mathbf{f}_4(A) \times \mathbf{f}_5(A) = \mathbf{f}_4(A) \times \mathbf{f}_5(A) \times \sum_e \mathbf{f}_2(E) \times \mathbf{f}_3(A, B, E).$$

Now the pointwise product inside the summation is computed, and the variable is summed out of the resulting matrix.

Notice that matrices are *not* multiplied until we need to sum out a variable from the accumulated product. At that point, we multiply just those matrices that include the variable to be summed out. Given functions for pointwise product and summing out, the variable elimination algorithm itself can be written quite simply, as shown in Figure 14.11.

### Variable ordering and variable relevance

The algorithm in Figure 14.11 includes an unspecified ORDER function to choose an ordering for the variables. Every choice of ordering yields a valid algorithm, but different orderings cause different intermediate factors to be generated during the calculation. For example, in the calculation shown previously, we eliminated  $A$  before  $E$ ; if we do it the other way, the calculation becomes

$$\mathbf{P}(B \mid j, m) = \alpha \mathbf{f}_1(B) \times \sum_a \mathbf{f}_4(A) \times \mathbf{f}_5(A) \times \sum_e \mathbf{f}_2(E) \times \mathbf{f}_3(A, B, E),$$

during which a new factor  $\mathbf{f}_6(A, B)$  will be generated.

In general, the time and space requirements of variable elimination are dominated by the size of the largest factor constructed during the operation of the algorithm. This in turn

```

function ELIMINATION-ASK( $X, \mathbf{e}, bn$ ) returns a distribution over  $X$ 
  inputs:  $X$ , the query variable
            $\mathbf{e}$ , observed values for variables  $\mathbf{E}$ 
            $bn$ , a Bayesian network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 

   $factors \leftarrow []$ 
  for each  $var$  in ORDER( $bn.VARS$ ) do
     $factors \leftarrow [MAKE-FACTOR(var, \mathbf{e}) | factors]$ 
    if  $var$  is a hidden variable then  $factors \leftarrow SUM-OUT(var, factors)$ 
  return NORMALIZE(POINTWISE-PRODUCT( $factors$ ))

```

**Figure 14.11** The variable elimination algorithm for inference in Bayesian networks.

is determined by the order of elimination of variables and by the structure of the network. It turns out to be intractable to determine the optimal ordering, but several good heuristics are available. One fairly effective method is a greedy one: eliminate whichever variable minimizes the size of the next factor to be constructed.

Let us consider one more query:  $\mathbf{P}(JohnCalls \mid Burglary = true)$ . As usual, the first step is to write out the nested summation:

$$\mathbf{P}(J \mid b) = \alpha P(b) \sum_e P(e) \sum_a P(a \mid b, e) \mathbf{P}(J \mid a) \sum_m P(m \mid a).$$

Evaluating this expression from right to left, we notice something interesting:  $\sum_m P(m \mid a)$  is equal to 1 by definition! Hence, there was no need to include it in the first place; the variable  $M$  is *irrelevant* to this query. Another way of saying this is that the result of the query  $P(JohnCalls \mid Burglary = true)$  is unchanged if we remove *MaryCalls* from the network altogether. In general, we can remove any leaf node that is not a query variable or an evidence variable. After its removal, there may be some more leaf nodes, and these too may be irrelevant. Continuing this process, we eventually find that *every variable that is not an ancestor of a query variable or evidence variable is irrelevant to the query*. A variable elimination algorithm can therefore remove all these variables before evaluating the query.

### 14.4.3 The complexity of exact inference

The complexity of exact inference in Bayesian networks depends strongly on the structure of the network. The burglary network of Figure 14.2 belongs to the family of networks in which there is at most one undirected path between any two nodes in the network. These are called **singly connected** networks or **polytrees**, and they have a particularly nice property: *The time and space complexity of exact inference in polytrees is linear in the size of the network*. Here, the size is defined as the number of CPT entries; if the number of parents of each node is bounded by a constant, then the complexity will also be linear in the number of nodes.

For **multiply connected** networks, such as that of Figure 14.12(a), variable elimination can have exponential time and space complexity in the worst case, even when the number of parents per node is bounded. This is not surprising when one considers that *because it*



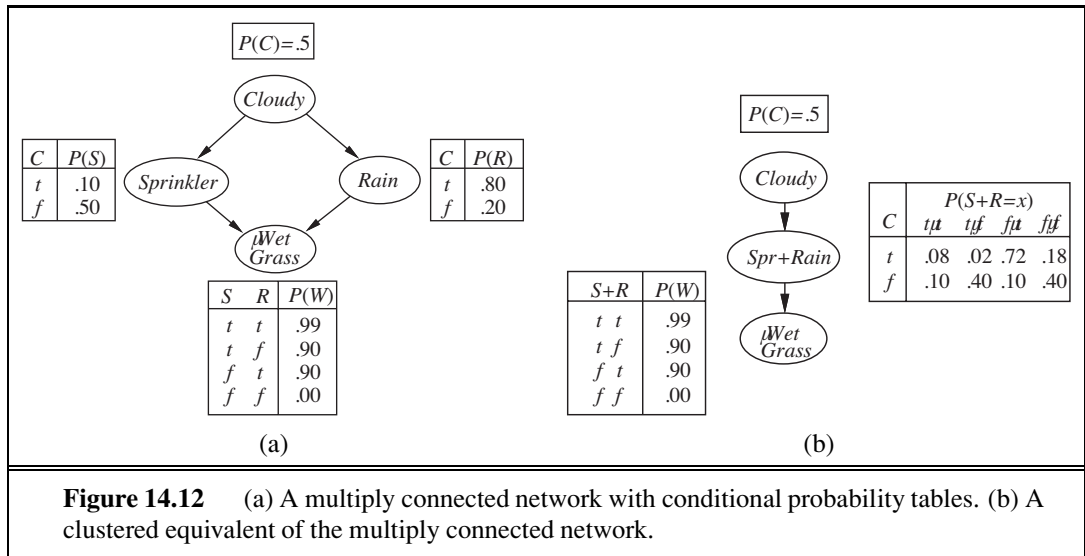
SINGLY CONNECTED

POLYTREE



MULTIPLY  
CONNECTED





includes inference in propositional logic as a special case, inference in Bayesian networks is NP-hard. In fact, it can be shown (Exercise 14.16) that the problem is as hard as that of computing the number of satisfying assignments for a propositional logic formula. This means that it is #P-hard (“number-P hard”)—that is, strictly harder than NP-complete problems.

There is a close connection between the complexity of Bayesian network inference and the complexity of constraint satisfaction problems (CSPs). As we discussed in Chapter 6, the difficulty of solving a discrete CSP is related to how “treelike” its constraint graph is. Measures such as **tree width**, which bound the complexity of solving a CSP, can also be applied directly to Bayesian networks. Moreover, the variable elimination algorithm can be generalized to solve CSPs as well as Bayesian networks.

#### 14.4.4 Clustering algorithms

The variable elimination algorithm is simple and efficient for answering individual queries. If we want to compute posterior probabilities for all the variables in a network, however, it can be less efficient. For example, in a polytree network, one would need to issue  $O(n)$  queries costing  $O(n)$  each, for a total of  $O(n^2)$  time. Using **clustering** algorithms (also known as **join tree** algorithms), the time can be reduced to  $O(n)$ . For this reason, these algorithms are widely used in commercial Bayesian network tools.

The basic idea of clustering is to join individual nodes of the network to form cluster nodes in such a way that the resulting network is a polytree. For example, the multiply connected network shown in Figure 14.12(a) can be converted into a polytree by combining the *Sprinkler* and *Rain* node into a cluster node called *Sprinkler+Rain*, as shown in Figure 14.12(b). The two Boolean nodes are replaced by a “meganode” that takes on four possible values: *tt*, *tf*, *ft*, and *ff*. The meganode has only one parent, the Boolean variable *Cloudy*, so there are two conditioning cases. Although this example doesn’t show it, the process of clustering often produces meganodes that share some variables.

Once the network is in polytree form, a special-purpose inference algorithm is required, because ordinary inference methods cannot handle meganodes that share variables with each other. Essentially, the algorithm is a form of constraint propagation (see Chapter 6) where the constraints ensure that neighboring meganodes agree on the posterior probability of any variables that they have in common. With careful bookkeeping, this algorithm is able to compute posterior probabilities for all the nonevidence nodes in the network in time *linear* in the size of the clustered network. However, the NP-hardness of the problem has not disappeared: if a network requires exponential time and space with variable elimination, then the CPTs in the clustered network will necessarily be exponentially large.

## 14.5 APPROXIMATE INFERENCE IN BAYESIAN NETWORKS

### MONTE CARLO

Given the intractability of exact inference in large, multiply connected networks, it is essential to consider approximate inference methods. This section describes randomized sampling algorithms, also called **Monte Carlo** algorithms, that provide approximate answers whose accuracy depends on the number of samples generated. Monte Carlo algorithms, of which simulated annealing (page 126) is an example, are used in many branches of science to estimate quantities that are difficult to calculate exactly. In this section, we are interested in sampling applied to the computation of posterior probabilities. We describe two families of algorithms: direct sampling and Markov chain sampling. Two other approaches—variational methods and loopy propagation—are mentioned in the notes at the end of the chapter.

### 14.5.1 Direct sampling methods

The primitive element in any sampling algorithm is the generation of samples from a known probability distribution. For example, an unbiased coin can be thought of as a random variable *Coin* with values  $\langle heads, tails \rangle$  and a prior distribution  $\mathbf{P}(Coin) = \langle 0.5, 0.5 \rangle$ . Sampling from this distribution is exactly like flipping the coin: with probability 0.5 it will return *heads*, and with probability 0.5 it will return *tails*. Given a source of random numbers uniformly distributed in the range  $[0, 1]$ , it is a simple matter to sample any distribution on a single variable, whether discrete or continuous. (See Exercise 14.17.)

The simplest kind of random sampling process for Bayesian networks generates events from a network that has no evidence associated with it. The idea is to sample each variable in turn, in topological order. The probability distribution from which the value is sampled is conditioned on the values already assigned to the variable's parents. This algorithm is shown in Figure 14.13. We can illustrate its operation on the network in Figure 14.12(a), assuming an ordering  $[Cloudy, Sprinkler, Rain, WetGrass]$ :

1. Sample from  $\mathbf{P}(Cloudy) = \langle 0.5, 0.5 \rangle$ , value is *true*.
2. Sample from  $\mathbf{P}(Sprinkler \mid Cloudy = true) = \langle 0.1, 0.9 \rangle$ , value is *false*.
3. Sample from  $\mathbf{P}(Rain \mid Cloudy = true) = \langle 0.8, 0.2 \rangle$ , value is *true*.
4. Sample from  $\mathbf{P}(WetGrass \mid Sprinkler = false, Rain = true) = \langle 0.9, 0.1 \rangle$ , value is *true*.

In this case, PRIOR-SAMPLE returns the event  $[true, false, true, true]$ .

```

function PRIOR-SAMPLE( $bn$ ) returns an event sampled from the prior specified by  $bn$ 
inputs:  $bn$ , a Bayesian network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 

 $\mathbf{x} \leftarrow$  an event with  $n$  elements
foreach variable  $X_i$  in  $X_1, \dots, X_n$  do
     $\mathbf{x}[i] \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$ 
return  $\mathbf{x}$ 

```

**Figure 14.13** A sampling algorithm that generates events from a Bayesian network. Each variable is sampled according to the conditional distribution given the values already sampled for the variable's parents.

It is easy to see that PRIOR-SAMPLE generates samples from the prior joint distribution specified by the network. First, let  $S_{PS}(x_1, \dots, x_n)$  be the probability that a specific event is generated by the PRIOR-SAMPLE algorithm. *Just looking at the sampling process*, we have

$$S_{PS}(x_1 \dots x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$

because each sampling step depends only on the parent values. This expression should look familiar, because it is also the probability of the event according to the Bayesian net's representation of the joint distribution, as stated in Equation (14.2). That is, we have

$$S_{PS}(x_1 \dots x_n) = P(x_1 \dots x_n) .$$

This simple fact makes it easy to answer questions by using samples.

In any sampling algorithm, the answers are computed by counting the actual samples generated. Suppose there are  $N$  total samples, and let  $N_{PS}(x_1, \dots, x_n)$  be the number of times the specific event  $x_1, \dots, x_n$  occurs in the set of samples. We expect this number, as a fraction of the total, to converge in the limit to its expected value according to the sampling probability:

$$\lim_{N \rightarrow \infty} \frac{N_{PS}(x_1, \dots, x_n)}{N} = S_{PS}(x_1, \dots, x_n) = P(x_1, \dots, x_n) . \quad (14.5)$$

For example, consider the event produced earlier:  $[true, false, true, true]$ . The sampling probability for this event is

$$S_{PS}(true, false, true, true) = 0.5 \times 0.9 \times 0.8 \times 0.9 = 0.324 .$$

Hence, in the limit of large  $N$ , we expect 32.4% of the samples to be of this event.

Whenever we use an approximate equality (" $\approx$ ") in what follows, we mean it in exactly this sense—that the estimated probability becomes exact in the large-sample limit. Such an estimate is called **consistent**. For example, one can produce a consistent estimate of the probability of any partially specified event  $x_1, \dots, x_m$ , where  $m \leq n$ , as follows:

$$P(x_1, \dots, x_m) \approx N_{PS}(x_1, \dots, x_m) / N . \quad (14.6)$$

That is, the probability of the event can be estimated as the fraction of all complete events generated by the sampling process that match the partially specified event. For example, if

we generate 1000 samples from the sprinkler network, and 511 of them have  $Rain = true$ , then the estimated probability of rain, written as  $\hat{P}(Rain = true)$ , is 0.511.

### Rejection sampling in Bayesian networks

REJECTION  
SAMPLING

**Rejection sampling** is a general method for producing samples from a hard-to-sample distribution given an easy-to-sample distribution. In its simplest form, it can be used to compute conditional probabilities—that is, to determine  $P(X | \mathbf{e})$ . The REJECTION-SAMPLING algorithm is shown in Figure 14.14. First, it generates samples from the prior distribution specified by the network. Then, it rejects all those that do not match the evidence. Finally, the estimate  $\hat{P}(X = x | \mathbf{e})$  is obtained by counting how often  $X = x$  occurs in the remaining samples.

Let  $\hat{\mathbf{P}}(X | \mathbf{e})$  be the estimated distribution that the algorithm returns. From the definition of the algorithm, we have

$$\hat{\mathbf{P}}(X | \mathbf{e}) = \alpha \mathbf{N}_{PS}(X, \mathbf{e}) = \frac{\mathbf{N}_{PS}(X, \mathbf{e})}{N_{PS}(\mathbf{e})}.$$

From Equation (14.6), this becomes

$$\hat{\mathbf{P}}(X | \mathbf{e}) \approx \frac{\mathbf{P}(X, \mathbf{e})}{P(\mathbf{e})} = \mathbf{P}(X | \mathbf{e}).$$

That is, rejection sampling produces a consistent estimate of the true probability.

Continuing with our example from Figure 14.12(a), let us assume that we wish to estimate  $\mathbf{P}(Rain | Sprinkler = true)$ , using 100 samples. Of the 100 that we generate, suppose that 73 have  $Sprinkler = false$  and are rejected, while 27 have  $Sprinkler = true$ ; of the 27, 8 have  $Rain = true$  and 19 have  $Rain = false$ . Hence,

$$\mathbf{P}(Rain | Sprinkler = true) \approx \text{NORMALIZE}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle.$$

The true answer is  $\langle 0.3, 0.7 \rangle$ . As more samples are collected, the estimate will converge to the true answer. The standard deviation of the error in each probability will be proportional to  $1/\sqrt{n}$ , where  $n$  is the number of samples used in the estimate.

The biggest problem with rejection sampling is that it rejects so many samples! The fraction of samples consistent with the evidence  $\mathbf{e}$  drops exponentially as the number of evidence variables grows, so the procedure is simply unusable for complex problems.

Notice that rejection sampling is very similar to the estimation of conditional probabilities directly from the real world. For example, to estimate  $\mathbf{P}(Rain | RedSkyAtNight = true)$ , one can simply count how often it rains after a red sky is observed the previous evening—ignoring those evenings when the sky is not red. (Here, the world itself plays the role of the sample-generation algorithm.) Obviously, this could take a long time if the sky is very seldom red, and that is the weakness of rejection sampling.

### Likelihood weighting

LIKELIHOOD  
WEIGHTING

**Likelihood weighting** avoids the inefficiency of rejection sampling by generating only events that are consistent with the evidence  $\mathbf{e}$ . It is a particular instance of the general statistical technique of **importance sampling**, tailored for inference in Bayesian networks. We begin by

IMPORTANCE  
SAMPLING



```

function REJECTION-SAMPLING( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $\mathbf{P}(X|\mathbf{e})$ 
  inputs:  $X$ , the query variable
            $\mathbf{e}$ , observed values for variables  $\mathbf{E}$ 
            $bn$ , a Bayesian network
            $N$ , the total number of samples to be generated
  local variables:  $\mathbf{N}$ , a vector of counts for each value of  $X$ , initially zero

  for  $j = 1$  to  $N$  do
     $\mathbf{x} \leftarrow \text{PRIOR-SAMPLE}(bn)$ 
    if  $\mathbf{x}$  is consistent with  $\mathbf{e}$  then
       $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $\mathbf{N}$ )

```

**Figure 14.14** The rejection-sampling algorithm for answering queries given evidence in a Bayesian network.

describing how the algorithm works; then we show that it works correctly—that is, generates consistent probability estimates.

LIKELIHOOD-WEIGHTING (see Figure 14.15) fixes the values for the evidence variables  $\mathbf{E}$  and samples only the nonevidence variables. This guarantees that each event generated is consistent with the evidence. Not all events are equal, however. Before tallying the counts in the distribution for the query variable, each event is weighted by the *likelihood* that the event accords to the evidence, as measured by the product of the conditional probabilities for each evidence variable, given its parents. Intuitively, events in which the actual evidence appears unlikely should be given less weight.

Let us apply the algorithm to the network shown in Figure 14.12(a), with the query  $\mathbf{P}(\text{Rain} \mid \text{Cloudy} = \text{true}, \text{WetGrass} = \text{true})$  and the ordering *Cloudy, Sprinkler, Rain, WetGrass*. (Any topological ordering will do.) The process goes as follows: First, the weight  $w$  is set to 1.0. Then an event is generated:

1. *Cloudy* is an evidence variable with value *true*. Therefore, we set

$$w \leftarrow w \times P(\text{Cloudy} = \text{true}) = 0.5 .$$

2. *Sprinkler* is not an evidence variable, so sample from  $\mathbf{P}(\text{Sprinkler} \mid \text{Cloudy} = \text{true}) = \langle 0.1, 0.9 \rangle$ ; suppose this returns *false*.
3. Similarly, sample from  $\mathbf{P}(\text{Rain} \mid \text{Cloudy} = \text{true}) = \langle 0.8, 0.2 \rangle$ ; suppose this returns *true*.
4. *WetGrass* is an evidence variable with value *true*. Therefore, we set

$$w \leftarrow w \times P(\text{WetGrass} = \text{true} \mid \text{Sprinkler} = \text{false}, \text{Rain} = \text{true}) = 0.45 .$$

Here WEIGHTED-SAMPLE returns the event  $[\text{true}, \text{false}, \text{true}, \text{true}]$  with weight 0.45, and this is tallied under *Rain = true*.

To understand why likelihood weighting works, we start by examining the sampling probability  $S_{WS}$  for WEIGHTED-SAMPLE. Remember that the evidence variables  $\mathbf{E}$  are fixed

---

```

function LIKELIHOOD-WEIGHTING( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $\mathbf{P}(X|\mathbf{e})$ 
  inputs:  $X$ , the query variable
            $\mathbf{e}$ , observed values for variables  $\mathbf{E}$ 
            $bn$ , a Bayesian network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 
            $N$ , the total number of samples to be generated
  local variables:  $\mathbf{W}$ , a vector of weighted counts for each value of  $X$ , initially zero

  for  $j = 1$  to  $N$  do
     $\mathbf{x}, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, \mathbf{e})$ 
     $\mathbf{W}[x] \leftarrow \mathbf{W}[x] + w$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $\mathbf{W}$ )

```

---

```

function WEIGHTED-SAMPLE( $bn, \mathbf{e}$ ) returns an event and a weight
   $w \leftarrow 1$ ;  $\mathbf{x} \leftarrow$  an event with  $n$  elements initialized from  $\mathbf{e}$ 
  foreach variable  $X_i$  in  $X_1, \dots, X_n$  do
    if  $X_i$  is an evidence variable with value  $x_i$  in  $\mathbf{e}$ 
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$ 
      else  $\mathbf{x}[i] \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$ 
  return  $\mathbf{x}, w$ 

```

---

**Figure 14.15** The likelihood-weighting algorithm for inference in Bayesian networks. In WEIGHTED-SAMPLE, each nonevidence variable is sampled according to the conditional distribution given the values already sampled for the variable's parents, while a weight is accumulated based on the likelihood for each evidence variable.

with values  $\mathbf{e}$ . We call the nonevidence variables  $\mathbf{Z}$  (including the query variable  $X$ ). The algorithm samples each variable in  $\mathbf{Z}$  given its parent values:

$$S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^l P(z_i \mid \text{parents}(Z_i)). \quad (14.7)$$

Notice that  $\text{Parents}(Z_i)$  can include both nonevidence variables and evidence variables. Unlike the prior distribution  $P(\mathbf{z})$ , the distribution  $S_{WS}$  pays some attention to the evidence: the sampled values for each  $Z_i$  will be influenced by evidence among  $Z_i$ 's ancestors. For example, when sampling *Sprinkler* the algorithm pays attention to the evidence *Cloudy* = *true* in its parent variable. On the other hand,  $S_{WS}$  pays less attention to the evidence than does the true posterior distribution  $P(\mathbf{z}|\mathbf{e})$ , because the sampled values for each  $Z_i$  ignore evidence among  $Z_i$ 's non-ancestors.<sup>5</sup> For example, when sampling *Sprinkler* and *Rain* the algorithm ignores the evidence in the child variable *WetGrass* = *true*; this means it will generate many samples with *Sprinkler* = *false* and *Rain* = *false* despite the fact that the evidence actually rules out this case.

---

<sup>5</sup> Ideally, we would like to use a sampling distribution equal to the true posterior  $P(\mathbf{z}|\mathbf{e})$ , to take all the evidence into account. This cannot be done efficiently, however. If it could, then we could approximate the desired probability to arbitrary accuracy with a polynomial number of samples. It can be shown that no such polynomial-time approximation scheme can exist.

The likelihood weight  $w$  makes up for the difference between the actual and desired sampling distributions. The weight for a given sample  $\mathbf{x}$ , composed from  $\mathbf{z}$  and  $\mathbf{e}$ , is the product of the likelihoods for each evidence variable given its parents (some or all of which may be among the  $Z_i$ s):

$$w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^m P(e_i \mid \text{parents}(E_i)) . \quad (14.8)$$

Multiplying Equations (14.7) and (14.8), we see that the *weighted* probability of a sample has the particularly convenient form

$$\begin{aligned} S_{WS}(\mathbf{z}, \mathbf{e}) w(\mathbf{z}, \mathbf{e}) &= \prod_{i=1}^l P(z_i \mid \text{parents}(Z_i)) \prod_{i=1}^m P(e_i \mid \text{parents}(E_i)) \\ &= P(\mathbf{z}, \mathbf{e}) \end{aligned} \quad (14.9)$$

because the two products cover all the variables in the network, allowing us to use Equation (14.2) for the joint probability.

Now it is easy to show that likelihood weighting estimates are consistent. For any particular value  $x$  of  $X$ , the estimated posterior probability can be calculated as follows:

$$\begin{aligned} \hat{P}(x \mid \mathbf{e}) &= \alpha \sum_{\mathbf{y}} N_{WS}(x, \mathbf{y}, \mathbf{e}) w(x, \mathbf{y}, \mathbf{e}) && \text{from LIKELIHOOD-WEIGHTING} \\ &\approx \alpha' \sum_{\mathbf{y}} S_{WS}(x, \mathbf{y}, \mathbf{e}) w(x, \mathbf{y}, \mathbf{e}) && \text{for large } N \\ &= \alpha' \sum_{\mathbf{y}} P(x, \mathbf{y}, \mathbf{e}) && \text{by Equation (14.9)} \\ &= \alpha' P(x, \mathbf{e}) = P(x \mid \mathbf{e}) . \end{aligned}$$

Hence, likelihood weighting returns consistent estimates.

Because likelihood weighting uses all the samples generated, it can be much more efficient than rejection sampling. It will, however, suffer a degradation in performance as the number of evidence variables increases. This is because most samples will have very low weights and hence the weighted estimate will be dominated by the tiny fraction of samples that accord more than an infinitesimal likelihood to the evidence. The problem is exacerbated if the evidence variables occur late in the variable ordering, because then the nonevidence variables will have no evidence in their parents and ancestors to guide the generation of samples. This means the samples will be simulations that bear little resemblance to the reality suggested by the evidence.

## 14.5.2 Inference by Markov chain simulation

**Markov chain Monte Carlo** (MCMC) algorithms work quite differently from rejection sampling and likelihood weighting. Instead of generating each sample from scratch, MCMC algorithms generate each sample by making a random change to the preceding sample. It is therefore helpful to think of an MCMC algorithm as being in a particular *current state* specifying a value for every variable and generating a *next state* by making random changes to the

## GIBBS SAMPLING

current state. (If this reminds you of simulated annealing from Chapter 4 or WALKSAT from Chapter 7, that is because both are members of the MCMC family.) Here we describe a particular form of MCMC called **Gibbs sampling**, which is especially well suited for Bayesian networks. (Other forms, some of them significantly more powerful, are discussed in the notes at the end of the chapter.) We will first describe what the algorithm does, then we will explain why it works.

### Gibbs sampling in Bayesian networks

The Gibbs sampling algorithm for Bayesian networks starts with an arbitrary state (with the evidence variables fixed at their observed values) and generates a next state by randomly sampling a value for one of the nonevidence variables  $X_i$ . The sampling for  $X_i$  is done *conditioned on the current values of the variables in the Markov blanket of  $X_i$* . (Recall from page 517 that the Markov blanket of a variable consists of its parents, children, and children's parents.) The algorithm therefore wanders randomly around the state space—the space of possible complete assignments—flipping one variable at a time, but keeping the evidence variables fixed.

Consider the query  $\mathbf{P}(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$  applied to the network in Figure 14.12(a). The evidence variables *Sprinkler* and *WetGrass* are fixed to their observed values and the nonevidence variables *Cloudy* and *Rain* are initialized randomly—let us say to *true* and *false* respectively. Thus, the initial state is  $[\text{true}, \text{true}, \text{false}, \text{true}]$ . Now the nonevidence variables are sampled repeatedly in an arbitrary order. For example:

1. *Cloudy* is sampled, given the current values of its Markov blanket variables: in this case, we sample from  $\mathbf{P}(\text{Cloudy} \mid \text{Sprinkler} = \text{true}, \text{Rain} = \text{false})$ . (Shortly, we will show how to calculate this distribution.) Suppose the result is *Cloudy* = *false*. Then the new current state is  $[\text{false}, \text{true}, \text{false}, \text{true}]$ .
2. *Rain* is sampled, given the current values of its Markov blanket variables: in this case, we sample from  $\mathbf{P}(\text{Rain} \mid \text{Cloudy} = \text{false}, \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$ . Suppose this yields *Rain* = *true*. The new current state is  $[\text{false}, \text{true}, \text{true}, \text{true}]$ .

Each state visited during this process is a sample that contributes to the estimate for the query variable *Rain*. If the process visits 20 states where *Rain* is true and 60 states where *Rain* is false, then the answer to the query is  $\text{NORMALIZE}(\langle 20, 60 \rangle) = \langle 0.25, 0.75 \rangle$ . The complete algorithm is shown in Figure 14.16.

### Why Gibbs sampling works

We will now show that Gibbs sampling returns consistent estimates for posterior probabilities. The material in this section is quite technical, but the basic claim is straightforward: *the sampling process settles into a “dynamic equilibrium” in which the long-run fraction of time spent in each state is exactly proportional to its posterior probability*. This remarkable property follows from the specific **transition probability** with which the process moves from one state to another, as defined by the conditional distribution given the Markov blanket of the variable being sampled.



TRANSITION  
PROBABILITY

```

function GIBBS-ASK( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $\mathbf{P}(X|\mathbf{e})$ 
  local variables:  $\mathbf{N}$ , a vector of counts for each value of  $X$ , initially zero
                    $\mathbf{Z}$ , the nonevidence variables in  $bn$ 
                    $\mathbf{x}$ , the current state of the network, initially copied from  $\mathbf{e}$ 

  initialize  $\mathbf{x}$  with random values for the variables in  $\mathbf{Z}$ 
  for  $j = 1$  to  $N$  do
    for each  $Z_i$  in  $\mathbf{Z}$  do
      set the value of  $Z_i$  in  $\mathbf{x}$  by sampling from  $\mathbf{P}(Z_i|mb(Z_i))$ 
       $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $\mathbf{N}$ )

```

**Figure 14.16** The Gibbs sampling algorithm for approximate inference in Bayesian networks; this version cycles through the variables, but choosing variables at random also works.

MARKOV CHAIN

Let  $q(\mathbf{x} \rightarrow \mathbf{x}')$  be the probability that the process makes a transition from state  $\mathbf{x}$  to state  $\mathbf{x}'$ . This transition probability defines what is called a **Markov chain** on the state space. (Markov chains also figure prominently in Chapters 15 and 17.) Now suppose that we run the Markov chain for  $t$  steps, and let  $\pi_t(\mathbf{x})$  be the probability that the system is in state  $\mathbf{x}$  at time  $t$ . Similarly, let  $\pi_{t+1}(\mathbf{x}')$  be the probability of being in state  $\mathbf{x}'$  at time  $t + 1$ . Given  $\pi_t(\mathbf{x})$ , we can calculate  $\pi_{t+1}(\mathbf{x}')$  by summing, for all states the system could be in at time  $t$ , the probability of being in that state times the probability of making the transition to  $\mathbf{x}'$ :

$$\pi_{t+1}(\mathbf{x}') = \sum_{\mathbf{x}} \pi_t(\mathbf{x}) q(\mathbf{x} \rightarrow \mathbf{x}').$$

STATIONARY DISTRIBUTION

We say that the chain has reached its **stationary distribution** if  $\pi_t = \pi_{t+1}$ . Let us call this stationary distribution  $\pi$ ; its defining equation is therefore

$$\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}) q(\mathbf{x} \rightarrow \mathbf{x}') \quad \text{for all } \mathbf{x}'. \quad (14.10)$$

ERGODIC

Provided the transition probability distribution  $q$  is **ergodic**—that is, every state is reachable from every other and there are no strictly periodic cycles—there is exactly one distribution  $\pi$  satisfying this equation for any given  $q$ .

Equation (14.10) can be read as saying that the expected “outflow” from each state (i.e., its current “population”) is equal to the expected “inflow” from all the states. One obvious way to satisfy this relationship is if the expected flow between any pair of states is the same in both directions; that is,

$$\pi(\mathbf{x}) q(\mathbf{x} \rightarrow \mathbf{x}') = \pi(\mathbf{x}') q(\mathbf{x}' \rightarrow \mathbf{x}) \quad \text{for all } \mathbf{x}, \mathbf{x}'. \quad (14.11)$$

DETAILED BALANCE

When these equations hold, we say that  $q(\mathbf{x} \rightarrow \mathbf{x}')$  is in **detailed balance** with  $\pi(\mathbf{x})$ .

We can show that detailed balance implies stationarity simply by summing over  $\mathbf{x}$  in Equation (14.11). We have

$$\sum_{\mathbf{x}} \pi(\mathbf{x}) q(\mathbf{x} \rightarrow \mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}') q(\mathbf{x}' \rightarrow \mathbf{x}) = \pi(\mathbf{x}') \sum_{\mathbf{x}} q(\mathbf{x}' \rightarrow \mathbf{x}) = \pi(\mathbf{x}')$$

where the last step follows because a transition from  $\mathbf{x}'$  is guaranteed to occur.

The transition probability  $q(\mathbf{x} \rightarrow \mathbf{x}')$  defined by the sampling step in GIBBS-ASK is actually a special case of the more general definition of Gibbs sampling, according to which each variable is sampled conditionally on the current values of *all* the other variables. We start by showing that this general definition of Gibbs sampling satisfies the detailed balance equation with a stationary distribution equal to  $P(\mathbf{x} | \mathbf{e})$ , (the true posterior distribution on the nonevidence variables). Then, we simply observe that, for Bayesian networks, sampling conditionally on all variables is equivalent to sampling conditionally on the variable's Markov blanket (see page 517).

To analyze the general Gibbs sampler, which samples each  $X_i$  in turn with a transition probability  $q_i$  that conditions on all the other variables, we define  $\bar{\mathbf{x}}_i$  to be these other variables (except the evidence variables); their values in the current state are  $\bar{\mathbf{x}}_i$ . If we sample a new value  $x'_i$  for  $X_i$  conditionally on all the other variables, including the evidence, we have

$$q_i(\mathbf{x} \rightarrow \mathbf{x}') = q_i((x_i, \bar{\mathbf{x}}_i) \rightarrow (x'_i, \bar{\mathbf{x}}_i)) = P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) .$$

Now we show that the transition probability for each step of the Gibbs sampler is in detailed balance with the true posterior:

$$\begin{aligned} \pi(\mathbf{x})q_i(\mathbf{x} \rightarrow \mathbf{x}') &= P(\mathbf{x} | \mathbf{e})P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) = P(x_i, \bar{\mathbf{x}}_i | \mathbf{e})P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) \\ &= P(x_i | \bar{\mathbf{x}}_i, \mathbf{e})P(\bar{\mathbf{x}}_i | \mathbf{e})P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) \quad (\text{using the chain rule on the first term}) \\ &= P(x_i | \bar{\mathbf{x}}_i, \mathbf{e})P(x'_i, \bar{\mathbf{x}}_i | \mathbf{e}) \quad (\text{using the chain rule backward}) \\ &= \pi(\mathbf{x}')q_i(\mathbf{x}' \rightarrow \mathbf{x}) . \end{aligned}$$

We can think of the loop “**for each**  $Z_i$  **in**  $\mathbf{Z}$  **do**” in Figure 14.16 as defining one large transition probability  $q$  that is the sequential composition  $q_1 \circ q_2 \circ \dots \circ q_n$  of the transition probabilities for the individual variables. It is easy to show (Exercise 14.19) that if each of  $q_i$  and  $q_j$  has  $\pi$  as its stationary distribution, then the sequential composition  $q_i \circ q_j$  does too; hence the transition probability  $q$  for the whole loop has  $P(\mathbf{x} | \mathbf{e})$  as its stationary distribution. Finally, unless the CPTs contain probabilities of 0 or 1—which can cause the state space to become disconnected—it is easy to see that  $q$  is ergodic. Hence, the samples generated by Gibbs sampling will eventually be drawn from the true posterior distribution.

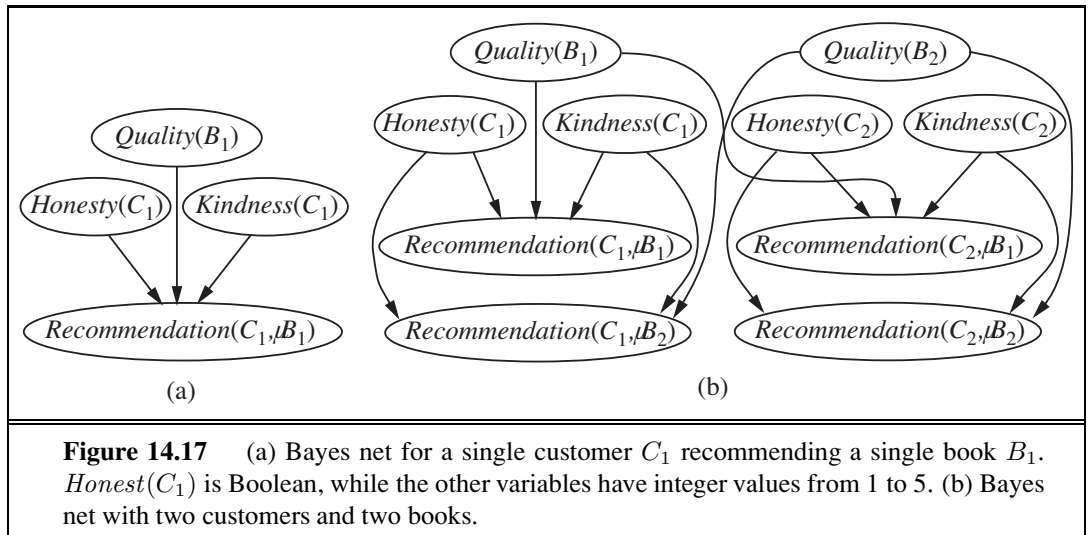
The final step is to show how to perform the general Gibbs sampling step—sampling  $X_i$  from  $\mathbf{P}(X_i | \bar{\mathbf{x}}_i, \mathbf{e})$ —in a Bayesian network. Recall from page 517 that a variable is independent of all other variables given its Markov blanket; hence,

$$P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) = P(x'_i | mb(X_i)) ,$$

where  $mb(X_i)$  denotes the values of the variables in  $X_i$ 's Markov blanket,  $MB(X_i)$ . As shown in Exercise 14.7, the probability of a variable given its Markov blanket is proportional to the probability of the variable given its parents times the probability of each child given its respective parents:

$$P(x'_i | mb(X_i)) = \alpha P(x'_i | parents(X_i)) \times \prod_{Y_j \in Children(X_i)} P(y_j | parents(Y_j)) . \quad (14.12)$$

Hence, to flip each variable  $X_i$  conditioned on its Markov blanket, the number of multiplications required is equal to the number of  $X_i$ 's children.



## 14.6 RELATIONAL AND FIRST-ORDER PROBABILITY MODELS



In Chapter 8, we explained the representational advantages possessed by first-order logic in comparison to propositional logic. First-order logic commits to the existence of objects and relations among them and can express facts about *some* or *all* of the objects in a domain. This often results in representations that are vastly more concise than the equivalent propositional descriptions. Now, Bayesian networks are essentially propositional: the set of random variables is fixed and finite, and each has a fixed domain of possible values. This fact limits the applicability of Bayesian networks. *If we can find a way to combine probability theory with the expressive power of first-order representations, we expect to be able to increase dramatically the range of problems that can be handled.*

For example, suppose that an online book retailer would like to provide overall evaluations of products based on recommendations received from its customers. The evaluation will take the form of a posterior distribution over the quality of the book, given the available evidence. The simplest solution to base the evaluation on the average recommendation, perhaps with a variance determined by the number of recommendations, but this fails to take into account the fact that some customers are kinder than others and some are less honest than others. Kind customers tend to give high recommendations even to fairly mediocre books, while dishonest customers give very high or very low recommendations for reasons other than quality—for example, they might work for a publisher.<sup>6</sup>

For a single customer  $C_1$ , recommending a single book  $B_1$ , the Bayes net might look like the one shown in Figure 14.17(a). (Just as in Section 9.1, expressions with parentheses such as  $Honest(C_1)$  are just fancy symbols—in this case, fancy names for random variables.)

<sup>6</sup> A game theorist would advise a dishonest customer to avoid detection by occasionally recommending a good book from a competitor. See Chapter 17.

With two customers and two books, the Bayes net looks like the one in Figure 14.17(b). For larger numbers of books and customers, it becomes completely impractical to specify the network by hand.

Fortunately, the network has a lot of repeated structure. Each  $Recommendation(c, b)$  variable has as its parents the variables  $Honest(c)$ ,  $Kindness(c)$ , and  $Quality(b)$ . Moreover, the CPTs for all the  $Recommendation(c, b)$  variables are identical, as are those for all the  $Honest(c)$  variables, and so on. The situation seems tailor-made for a first-order language. We would like to say something like

$$Recommendation(c, b) \sim RecCPT(Honest(c), Kindness(c), Quality(b))$$

with the intended meaning that a customer's recommendation for a book depends on the customer's honesty and kindness and the book's quality according to some fixed CPT. This section develops a language that lets us say exactly this, and a lot more besides.

### 14.6.1 Possible worlds

Recall from Chapter 13 that a probability model defines a set  $\Omega$  of possible worlds with a probability  $P(\omega)$  for each world  $\omega$ . For Bayesian networks, the possible worlds are assignments of values to variables; for the Boolean case in particular, the possible worlds are identical to those of propositional logic. For a first-order probability model, then, it seems we need the possible worlds to be those of first-order logic—that is, a set of objects with relations among them and an interpretation that maps constant symbols to objects, predicate symbols to relations, and function symbols to functions on those objects. (See Section 8.2.) The model also needs to define a probability for each such possible world, just as a Bayesian network defines a probability for each assignment of values to variables.

Let us suppose, for a moment, that we have figured out how to do this. Then, as usual (see page 485), we can obtain the probability of any first-order logical sentence  $\phi$  as a sum over the possible worlds where it is true:

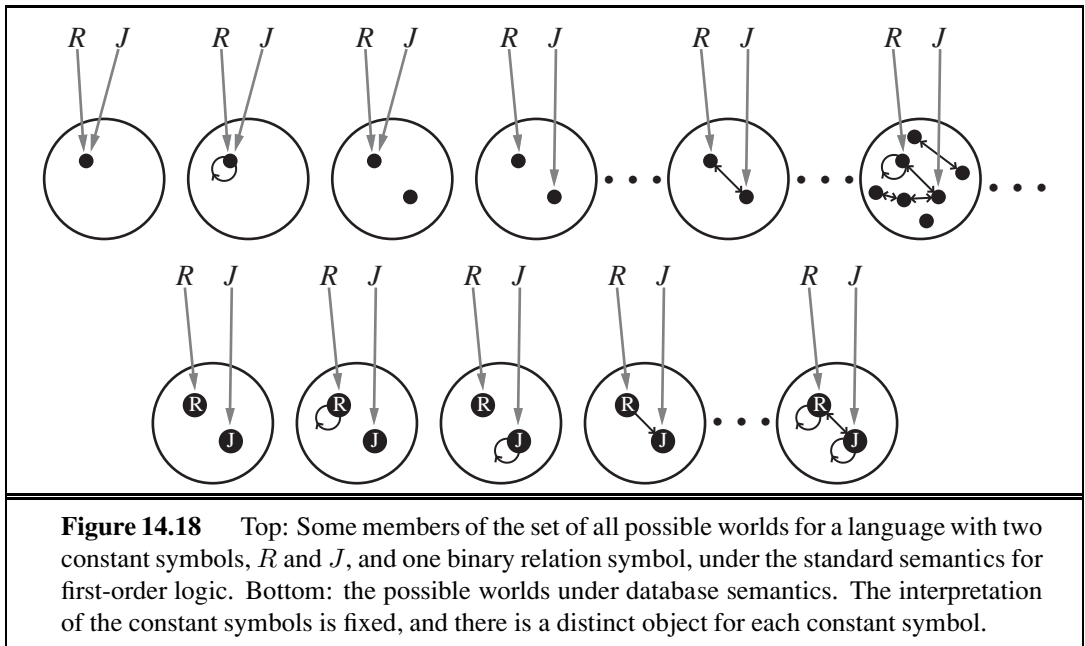
$$P(\phi) = \sum_{\omega: \phi \text{ is true in } \omega} P(\omega) . \quad (14.13)$$

Conditional probabilities  $P(\phi | \mathbf{e})$  can be obtained similarly, so we can, in principle, ask any question we want of our model—e.g., “Which books are most likely to be recommended highly by dishonest customers?”—and get an answer. So far, so good.

There is, however, a problem: the set of first-order models is infinite. We saw this explicitly in Figure 8.4 on page 293, which we show again in Figure 14.18 (top). This means that (1) the summation in Equation (14.13) could be infeasible, and (2) specifying a complete, consistent distribution over an infinite set of worlds could be very difficult.

Section 14.6.2 explores one approach to dealing with this problem. The idea is to borrow not from the standard semantics of first-order logic but from the **database semantics** defined in Section 8.2.8 (page 299). The database semantics makes the **unique names assumption**—here, we adopt it for the constant symbols. It also assumes **domain closure**—there are no more objects than those that are named. We can then guarantee a finite set of possible worlds by making the set of objects in each world be exactly the set of constant





symbols that are used; as shown in Figure 14.18 (bottom), there is no uncertainty about the mapping from symbols to objects or about the objects that exist. We will call models defined in this way **relational probability models**, or RPMs.<sup>7</sup> The most significant difference between the semantics of RPMs and the database semantics introduced in Section 8.2.8 is that RPMs do not make the closed-world assumption—obviously, assuming that every unknown fact is false doesn’t make sense in a probabilistic reasoning system!

When the underlying assumptions of database semantics fail to hold, RPMs won’t work well. For example, a book retailer might use an ISBN (International Standard Book Number) as a constant symbol to name each book, even though a given “logical” book (e.g., “Gone With the Wind”) may have several ISBNs. It would make sense to aggregate recommendations across multiple ISBNs, but the retailer may not know for sure which ISBNs are really the same book. (Note that we are not reifying the *individual copies* of the book, which might be necessary for used-book sales, car sales, and so on.) Worse still, each customer is identified by a login ID, but a dishonest customer may have thousands of IDs! In the computer security field, these multiple IDs are called **sibyls** and their use to confound a reputation system is called a **sibyl attack**. Thus, even a simple application in a relatively well-defined, online domain involves both **existence uncertainty** (what are the real books and customers underlying the observed data) and **identity uncertainty** (which symbol really refer to the same object). We need to bite the bullet and define probability models based on the standard semantics of first-order logic, for which the possible worlds vary in the objects they contain and in the mappings from symbols to objects. Section 14.6.3 shows how to do this.

<sup>7</sup> The name *relational probability model* was given by Pfeffer (2000) to a slightly different representation, but the underlying ideas are the same.

### 14.6.2 Relational probability models

TYPE SIGNATURE

Like first-order logic, RPMs have constant, function, and predicate symbols. (It turns out to be easier to view predicates as functions that return *true* or *false*.) We will also assume a **type signature** for each function, that is, a specification of the type of each argument and the function's value. If the type of each object is known, many spurious possible worlds are eliminated by this mechanism. For the book-recommendation domain, the types are *Customer* and *Book*, and the type signatures for the functions and predicates are as follows:

$$\text{Honest} : \text{Customer} \rightarrow \{\text{true}, \text{false}\} \quad \text{Kindness} : \text{Customer} \rightarrow \{1, 2, 3, 4, 5\}$$

$$\text{Quality} : \text{Book} \rightarrow \{1, 2, 3, 4, 5\}$$

$$\text{Recommendation} : \text{Customer} \times \text{Book} \rightarrow \{1, 2, 3, 4, 5\}$$

The constant symbols will be whatever customer and book names appear in the retailer's data set. In the example given earlier (Figure 14.17(b)), these were  $C_1$ ,  $C_2$  and  $B_1$ ,  $B_2$ .

Given the constants and their types, together with the functions and their type signatures, the random variables of the RPM are obtained by instantiating each function with each possible combination of objects:  $\text{Honest}(C_1)$ ,  $\text{Quality}(B_2)$ ,  $\text{Recommendation}(C_1, B_2)$ , and so on. These are exactly the variables appearing in Figure 14.17(b). Because each type has only finitely many instances, the number of basic random variables is also finite.

To complete the RPM, we have to write the dependencies that govern these random variables. There is one dependency statement for each function, where each argument of the function is a logical variable (i.e., a variable that ranges over objects, as in first-order logic):

$$\text{Honest}(c) \sim \langle 0.99, 0.01 \rangle$$

$$\text{Kindness}(c) \sim \langle 0.1, 0.1, 0.2, 0.3, 0.3 \rangle$$

$$\text{Quality}(b) \sim \langle 0.05, 0.2, 0.4, 0.2, 0.15 \rangle$$

$$\text{Recommendation}(c, b) \sim \text{RecCPT}(\text{Honest}(c), \text{Kindness}(c), \text{Quality}(b))$$

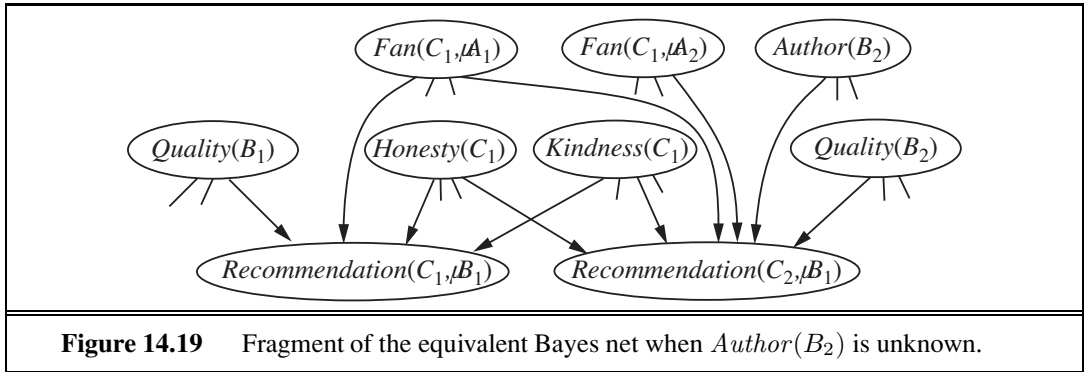
where *RecCPT* is a separately defined conditional distribution with  $2 \times 5 \times 5 = 50$  rows, each with 5 entries. The semantics of the RPM can be obtained by instantiating these dependencies for all known constants, giving a Bayesian network (as in Figure 14.17(b)) that defines a joint distribution over the RPM's random variables.<sup>8</sup>

CONTEXT-SPECIFIC  
INDEPENDENCE

We can refine the model by introducing a **context-specific independence** to reflect the fact that dishonest customers ignore quality when giving a recommendation; moreover, kindness plays no role in their decisions. A context-specific independence allows a variable to be independent of some of its parents given certain values of others; thus,  $\text{Recommendation}(c, b)$  is independent of  $\text{Kindness}(c)$  and  $\text{Quality}(b)$  when  $\text{Honest}(c) = \text{false}$ :

$$\begin{aligned} \text{Recommendation}(c, b) \sim & \text{if } \text{Honest}(c) \text{ then} \\ & \text{HonestRecCPT}(\text{Kindness}(c), \text{Quality}(b)) \\ & \text{else } \langle 0.4, 0.1, 0.0, 0.1, 0.4 \rangle . \end{aligned}$$

<sup>8</sup> Some technical conditions must be observed to guarantee that the RPM defines a proper distribution. First, the dependencies must be *acyclic*, otherwise the resulting Bayesian network will have cycles and will not define a proper distribution. Second, the dependencies must be *well-founded*, that is, there can be no infinite ancestor chains, such as might arise from recursive dependencies. Under some circumstances (see Exercise 14.6), a fixed-point calculation yields a well-defined probability model for a recursive RPM.



This kind of dependency may look like an ordinary if–then–else statement on a programming language, but there is a key difference: the inference engine *doesn't necessarily know the value of the conditional test*!

We can elaborate this model in endless ways to make it more realistic. For example, suppose that an honest customer who is a fan of a book's author always gives the book a 5, regardless of quality:

$$\begin{aligned}
 Recommendation(c, b) \sim & \text{ if } Honest(c) \text{ then} \\
 & \text{ if } Fan(c, Author(b)) \text{ then Exactly}(5) \\
 & \text{ else } HonestRecCPT(Kindness(c), Quality(b)) \\
 & \text{ else } \langle 0.4, 0.1, 0.0, 0.1, 0.4 \rangle
 \end{aligned}$$

Again, the conditional test  $Fan(c, Author(b))$  is unknown, but if a customer gives only 5s to a particular author's books and is not otherwise especially kind, then the posterior probability that the customer is a fan of that author will be high. Furthermore, the posterior distribution will tend to discount the customer's 5s in evaluating the quality of that author's books.

In the preceding example, we implicitly assumed that the value of  $Author(b)$  is known for every  $b$ , but this may not be the case. How can the system reason about whether, say,  $C_1$  is a fan of  $Author(B_2)$  when  $Author(B_2)$  is unknown? The answer is that the system may have to reason about *all possible authors*. Suppose (to keep things simple) that there are just two authors,  $A_1$  and  $A_2$ . Then  $Author(B_2)$  is a random variable with two possible values,  $A_1$  and  $A_2$ , and it is a parent of  $Recommendation(C_1, B_2)$ . The variables  $Fan(C_1, A_1)$  and  $Fan(C_1, A_2)$  are parents too. The conditional distribution for  $Recommendation(C_1, B_2)$  is then essentially a **multiplexer** in which the  $Author(B_2)$  parent acts as a selector to choose which of  $Fan(C_1, A_1)$  and  $Fan(C_1, A_2)$  actually gets to influence the recommendation. A fragment of the equivalent Bayes net is shown in Figure 14.19. Uncertainty in the value of  $Author(B_2)$ , which affects the dependency structure of the network, is an instance of **relational uncertainty**.

In case you are wondering how the system can possibly work out who the author of  $B_2$  is: consider the possibility that three other customers are fans of  $A_1$  (and have no other favorite authors in common) and all three have given  $B_2$  a 5, even though most other customers find it quite dismal. In that case, it is extremely likely that  $A_1$  is the author of  $B_2$ .

MULTIPLEXER

RELATIONAL  
UNCERTAINTY

The emergence of sophisticated reasoning like this from an RPM model of just a few lines is an intriguing example of how probabilistic influences spread through the web of interconnections among objects in the model. As more dependencies and more objects are added, the picture conveyed by the posterior distribution often becomes clearer and clearer.

UNROLLING

The next question is how to do inference in RPMs. One approach is to collect the evidence and query and the constant symbols therein, construct the equivalent Bayes net, and apply any of the inference methods discussed in this chapter. This technique is called **unrolling**. The obvious drawback is that the resulting Bayes net may be very large. Furthermore, if there are many candidate objects for an unknown relation or function—for example, the unknown author of  $B_2$ —then some variables in the network may have many parents.

Fortunately, much can be done to improve on generic inference algorithms. First, the presence of repeated substructure in the unrolled Bayes net means that many of the factors constructed during variable elimination (and similar kinds of tables constructed by clustering algorithms) will be identical; effective caching schemes have yielded speedups of three orders of magnitude for large networks. Second, inference methods developed to take advantage of context-specific independence in Bayes nets find many applications in RPMs. Third, MCMC inference algorithms have some interesting properties when applied to RPMs with relational uncertainty. MCMC works by sampling complete possible worlds, so in each state the relational structure is completely known. In the example given earlier, each MCMC state would specify the value of  $Author(B_2)$ , and so the other potential authors are no longer parents of the recommendation nodes for  $B_2$ . For MCMC, then, relational uncertainty causes no increase in network complexity; instead, the MCMC process includes transitions that change the relational structure, and hence the dependency structure, of the unrolled network.

All of the methods just described assume that the RPM has to be partially or completely unrolled into a Bayesian network. This is exactly analogous to the method of **proposition-alization** for first-order logical inference. (See page 322.) Resolution theorem-provers and logic programming systems avoid propositionalizing by instantiating the logical variables only as needed to make the inference go through; that is, they *lift* the inference process above the level of ground propositional sentences and make each lifted step do the work of many ground steps. The same idea applied in probabilistic inference. For example, in the variable elimination algorithm, a lifted factor can represent an entire set of ground factors that assign probabilities to random variables in the RPM, where those random variables differ only in the constant symbols used to construct them. The details of this method are beyond the scope of this book, but references are given at the end of the chapter.

### 14.6.3 Open-universe probability models

We argued earlier that database semantics was appropriate for situations in which we know exactly the set of relevant objects that exist and can identify them unambiguously. (In particular, all observations about an object are correctly associated with the constant symbol that names it.) In many real-world settings, however, these assumptions are simply untenable. We gave the examples of multiple ISBNs and sibyl attacks in the book-recommendation domain (to which we will return in a moment), but the phenomenon is far more pervasive:

- A vision system doesn't know what exists, if anything, around the next corner, and may not know if the object it sees now is the same one it saw a few minutes ago.
- A text-understanding system does not know in advance the entities that will be featured in a text, and must reason about whether phrases such as “Mary,” “Dr. Smith,” “she,” “his cardiologist,” “his mother,” and so on refer to the same object.
- An intelligence analyst hunting for spies never knows how many spies there really are and can only guess whether various pseudonyms, phone numbers, and sightings belong to the same individual.

In fact, a major part of human cognition seems to require learning what objects exist and being able to connect observations—which almost never come with unique IDs attached—to hypothesized objects in the world.

OPEN UNIVERSE

For these reasons, we need to be able to write so-called **open-universe** probability models or OUPMs based on the standard semantics of first-order logic, as illustrated at the top of Figure 14.18. A language for OUPMs provides a way of writing such models easily while guaranteeing a unique, consistent probability distribution over the infinite space of possible worlds.

The basic idea is to understand how ordinary Bayesian networks and RPMs manage to define a unique probability model and to transfer that insight to the first-order setting. In essence, a Bayes net *generates* each possible world, event by event, in the topological order defined by the network structure, where each event is an assignment of a value to a variable. An RPM extends this to entire sets of events, defined by the possible instantiations of the logical variables in a given predicate or function. OUPMs go further by allowing generative steps that *add objects* to the possible world under construction, where the number and type of objects may depend on the objects that are already in that world. That is, the event being generated is not the assignment of a value to a variable, but the very *existence* of objects.

One way to do this in OUPMs is to add statements that define conditional distributions over the numbers of objects of various kinds. For example, in the book-recommendation domain, we might want to distinguish between *customers* (real people) and their *login IDs*. Suppose we expect somewhere between 100 and 10,000 distinct customers (whom we cannot observe directly). We can express this as a prior log-normal distribution<sup>9</sup> as follows:

$$\# \text{ Customer} \sim \text{LogNormal}[6.9, 2.3^2]() .$$

We expect honest customers to have just one ID, whereas dishonest customers might have anywhere between 10 and 1000 IDs:

$$\# \text{ LoginID}(\text{Owner} = c) \sim \begin{array}{ll} \text{if } \text{Honest}(c) \text{ then } \text{Exactly}(1) \\ \text{else } \text{LogNormal}[6.9, 2.3^2]() . \end{array}$$

ORIGIN FUNCTION

This statement defines the number of login IDs for a given owner, who is a customer. The *Owner* function is called an **origin function** because it says where each generated object came from. In the formal semantics of BLOG (as distinct from first-order logic), the domain elements in each possible world are actually generation histories (e.g., “the fourth login ID of the seventh customer”) rather than simple tokens.

<sup>9</sup> A distribution  $\text{LogNormal}[\mu, \sigma^2](x)$  is equivalent to a distribution  $N[\mu, \sigma^2](x)$  over  $\log_e(x)$ .

Subject to technical conditions of acyclicity and well-foundedness similar to those for RPMs, open-universe models of this kind define a unique distribution over possible worlds. Furthermore, there exist inference algorithms such that, for every such well-defined model and every first-order query, the answer returned approaches the true posterior arbitrarily closely in the limit. There are some tricky issues involved in designing these algorithms. For example, an MCMC algorithm cannot sample directly in the space of possible worlds when the size of those worlds is unbounded; instead, it samples finite, partial worlds, relying on the fact that only finitely many objects can be relevant to the query in distinct ways. Moreover, transitions must allow for merging two objects into one or splitting one into two. (Details are given in the references at the end of the chapter.) Despite these complications, the basic principle established in Equation (14.13) still holds: the probability of any sentence is well defined and can be calculated.

Research in this area is still at an early stage, but already it is becoming clear that first-order probabilistic reasoning yields a tremendous increase in the effectiveness of AI systems at handling uncertain information. Potential applications include those mentioned above—computer vision, text understanding, and intelligence analysis—as well as many other kinds of sensor interpretation.

---

## 14.7 OTHER APPROACHES TO UNCERTAIN REASONING

---

Other sciences (e.g., physics, genetics, and economics) have long favored probability as a model for uncertainty. In 1819, Pierre Laplace said, “Probability theory is nothing but common sense reduced to calculation.” In 1850, James Maxwell said, “The true logic for this world is the calculus of Probabilities, which takes account of the magnitude of the probability which is, or ought to be, in a reasonable man’s mind.”

Given this long tradition, it is perhaps surprising that AI has considered many alternatives to probability. The earliest expert systems of the 1970s ignored uncertainty and used strict logical reasoning, but it soon became clear that this was impractical for most real-world domains. The next generation of expert systems (especially in medical domains) used probabilistic techniques. Initial results were promising, but they did not scale up because of the exponential number of probabilities required in the full joint distribution. (Efficient Bayesian network algorithms were unknown then.) As a result, probabilistic approaches fell out of favor from roughly 1975 to 1988, and a variety of alternatives to probability were tried for a variety of reasons:

- One common view is that probability theory is essentially numerical, whereas human judgmental reasoning is more “qualitative.” Certainly, we are not consciously aware of doing numerical calculations of degrees of belief. (Neither are we aware of doing unification, yet we seem to be capable of some kind of logical reasoning.) It might be that we have some kind of numerical degrees of belief encoded directly in strengths of connections and activations in our neurons. In that case, the difficulty of conscious access to those strengths is not surprising. One should also note that qualitative reason-

ing mechanisms can be built directly on top of probability theory, so the “no numbers” argument against probability has little force. Nonetheless, some qualitative schemes have a good deal of appeal in their own right. One of the best studied is **default reasoning**, which treats conclusions not as “believed to a certain degree,” but as “believed until a better reason is found to believe something else.” Default reasoning is covered in Chapter 12.

- **Rule-based** approaches to uncertainty have also been tried. Such approaches hope to build on the success of logical rule-based systems, but add a sort of “fudge factor” to each rule to accommodate uncertainty. These methods were developed in the mid-1970s and formed the basis for a large number of expert systems in medicine and other areas.
- One area that we have not addressed so far is the question of **ignorance**, as opposed to uncertainty. Consider the flipping of a coin. If we know that the coin is fair, then a probability of 0.5 for heads is reasonable. If we know that the coin is biased, but we do not know which way, then 0.5 for heads is again reasonable. Obviously, the two cases are different, yet the outcome probability seems not to distinguish them. The **Dempster–Shafer theory** uses **interval-valued** degrees of belief to represent an agent’s knowledge of the probability of a proposition.
- Probability makes the same ontological commitment as logic: that propositions are true or false in the world, even if the agent is uncertain as to which is the case. Researchers in **fuzzy logic** have proposed an ontology that allows **vagueness**: that a proposition can be “sort of” true. Vagueness and uncertainty are in fact orthogonal issues.

The next three subsections treat some of these approaches in slightly more depth. We will not provide detailed technical material, but we cite references for further study.

### 14.7.1 Rule-based methods for uncertain reasoning

Rule-based systems emerged from early work on practical and intuitive systems for logical inference. Logical systems in general, and logical rule-based systems in particular, have three desirable properties:

LOCALITY

- **Locality**: In logical systems, whenever we have a rule of the form  $A \Rightarrow B$ , we can conclude  $B$ , given evidence  $A$ , *without worrying about any other rules*. In probabilistic systems, we need to consider *all* the evidence.

DETACHMENT

- **Detachment**: Once a logical proof is found for a proposition  $B$ , the proposition can be used regardless of how it was derived. That is, it can be **detached** from its justification. In dealing with probabilities, on the other hand, the source of the evidence for a belief is important for subsequent reasoning.

TRUTH-FUNCTIONALITY

- **Truth-functionality**: In logic, the truth of complex sentences can be computed from the truth of the components. Probability combination does not work this way, except under strong global independence assumptions.

There have been several attempts to devise uncertain reasoning schemes that retain these advantages. The idea is to attach degrees of belief to propositions and rules and to devise purely local schemes for combining and propagating those degrees of belief. The schemes



are also truth-functional; for example, the degree of belief in  $A \vee B$  is a function of the belief in  $A$  and the belief in  $B$ .

The bad news for rule-based systems is that the properties of *locality*, *detachment*, and *truth-functionality* are simply not appropriate for uncertain reasoning. Let us look at truth-functionality first. Let  $H_1$  be the event that a fair coin flip comes up heads, let  $T_1$  be the event that the coin comes up tails on that same flip, and let  $H_2$  be the event that the coin comes up heads on a second flip. Clearly, all three events have the same probability, 0.5, and so a truth-functional system must assign the same belief to the disjunction of any two of them. But we can see that the probability of the disjunction depends on the events themselves and not just on their probabilities:

$P(A)$	$P(B)$	$P(A \vee B)$
$P(H_1) = 0.5$	$P(H_1) = 0.5$	$P(H_1 \vee H_1) = 0.50$
	$P(T_1) = 0.5$	$P(H_1 \vee T_1) = 1.00$
	$P(H_2) = 0.5$	$P(H_1 \vee H_2) = 0.75$

It gets worse when we chain evidence together. Truth-functional systems have **rules** of the form  $A \mapsto B$  that allow us to compute the belief in  $B$  as a function of the belief in the rule and the belief in  $A$ . Both forward- and backward-chaining systems can be devised. The belief in the rule is assumed to be constant and is usually specified by the knowledge engineer—for example, as  $A \mapsto_{0.9} B$ .

Consider the wet-grass situation from Figure 14.12(a) (page 529). If we wanted to be able to do both causal and diagnostic reasoning, we would need the two rules

$$Rain \mapsto WetGrass \quad \text{and} \quad WetGrass \mapsto Rain .$$

These two rules form a feedback loop: evidence for *Rain* increases the belief in *WetGrass*, which in turn increases the belief in *Rain* even more. Clearly, uncertain reasoning systems need to keep track of the paths along which evidence is propagated.

Intercausal reasoning (or explaining away) is also tricky. Consider what happens when we have the two rules

$$Sprinkler \mapsto WetGrass \quad \text{and} \quad WetGrass \mapsto Rain .$$

Suppose we see that the sprinkler is on. Chaining forward through our rules, this increases the belief that the grass will be wet, which in turn increases the belief that it is raining. But this is ridiculous: the fact that the sprinkler is on explains away the wet grass and should *reduce* the belief in rain. A truth-functional system acts as if it also believes  $Sprinkler \mapsto Rain$ .

Given these difficulties, how can truth-functional systems be made useful in practice? The answer lies in restricting the task and in carefully engineering the rule base so that undesirable interactions do not occur. The most famous example of a truth-functional system for uncertain reasoning is the **certainty factors** model, which was developed for the MYCIN medical diagnosis program and was widely used in expert systems of the late 1970s and 1980s. Almost all uses of certainty factors involved rule sets that were either purely diagnostic (as in MYCIN) or purely causal. Furthermore, evidence was entered only at the “roots” of the rule set, and most rule sets were singly connected. Heckerman (1986) has shown that,



under these circumstances, a minor variation on certainty-factor inference was exactly equivalent to Bayesian inference on polytrees. In other circumstances, certainty factors could yield disastrously incorrect degrees of belief through overcounting of evidence. As rule sets became larger, undesirable interactions between rules became more common, and practitioners found that the certainty factors of many other rules had to be “tweaked” when new rules were added. For these reasons, Bayesian networks have largely supplanted rule-based methods for uncertain reasoning.

### 14.7.2 Representing ignorance: Dempster–Shafer theory

DEMPSTER-SHAFFER  
THEORY

The **Dempster–Shafer theory** is designed to deal with the distinction between **uncertainty** and **ignorance**. Rather than computing the probability of a proposition, it computes the probability that the evidence supports the proposition. This measure of belief is called a **belief function**, written  $Bel(X)$ .

BELIEF FUNCTION

We return to coin flipping for an example of belief functions. Suppose you pick a coin from a magician’s pocket. Given that the coin might or might not be fair, what belief should you ascribe to the event that it comes up heads? Dempster–Shafer theory says that because you have no evidence either way, you have to say that the belief  $Bel(Heads) = 0$  and also that  $Bel(\neg Heads) = 0$ . This makes Dempster–Shafer reasoning systems skeptical in a way that has some intuitive appeal. Now suppose you have an expert at your disposal who testifies with 90% certainty that the coin is fair (i.e., he is 90% sure that  $P(Heads) = 0.5$ ). Then Dempster–Shafer theory gives  $Bel(Heads) = 0.9 \times 0.5 = 0.45$  and likewise  $Bel(\neg Heads) = 0.45$ . There is still a 10 percentage point “gap” that is not accounted for by the evidence.

MASS

The mathematical underpinnings of Dempster–Shafer theory have a similar flavor to those of probability theory; the main difference is that, instead of assigning probabilities to possible worlds, the theory assigns **masses** to *sets* of possible world, that is, to events. The masses still must add to 1 over all possible events.  $Bel(A)$  is defined to be the sum of masses for all events that are subsets of (i.e., that entail)  $A$ , including  $A$  itself. With this definition,  $Bel(A)$  and  $Bel(\neg A)$  sum to *at most* 1, and the gap—the interval between  $Bel(A)$  and  $1 - Bel(\neg A)$ —is often interpreted as bounding the probability of  $A$ .

As with default reasoning, there is a problem in connecting beliefs to actions. Whenever there is a gap in the beliefs, then a decision problem can be defined such that a Dempster–Shafer system is unable to make a decision. In fact, the notion of utility in the Dempster–Shafer model is not yet well understood because the meanings of masses and beliefs themselves have yet to be understood. Pearl (1988) has argued that  $Bel(A)$  should be interpreted not as a degree of belief in  $A$  but as the probability assigned to all the possible worlds (now interpreted as logical theories) in which  $A$  is *provable*. While there are cases in which this quantity might be of interest, it is not the same as the probability that  $A$  is true.

A Bayesian analysis of the coin-flipping example would suggest that no new formalism is necessary to handle such cases. The model would have two variables: the *Bias* of the coin (a number between 0 and 1, where 0 is a coin that always shows tails and 1 a coin that always shows heads) and the outcome of the next *Flip*. The prior probability distribution for *Bias*

would reflect our beliefs based on the source of the coin (the magician's pocket): some small probability that it is fair and some probability that it is heavily biased toward heads or tails. The conditional distribution  $\mathbf{P}(\text{Flip} \mid \text{Bias})$  simply defines how the bias operates. If  $\mathbf{P}(\text{Bias})$  is symmetric about 0.5, then our prior probability for the flip is

$$P(\text{Flip} = \text{heads}) = \int_0^1 P(\text{Bias} = x)P(\text{Flip} = \text{heads} \mid \text{Bias} = x) dx = 0.5 .$$

This is the same prediction as if we believe strongly that the coin is fair, but that does *not* mean that probability theory treats the two situations identically. The difference arises *after* the flips in computing the posterior distribution for *Bias*. If the coin came from a bank, then seeing it come up heads three times running would have almost no effect on our strong prior belief in its fairness; but if the coin comes from the magician's pocket, the same evidence will lead to a stronger posterior belief that the coin is biased toward heads. Thus, a Bayesian approach expresses our “ignorance” in terms of how our beliefs would change in the face of future information gathering.

### 14.7.3 Representing vagueness: Fuzzy sets and fuzzy logic

#### FUZZY SET THEORY



**Fuzzy set theory** is a means of specifying how well an object satisfies a vague description. For example, consider the proposition “Nate is tall.” Is this true if Nate is 5' 10"? Most people would hesitate to answer “true” or “false,” preferring to say, “sort of.” Note that this is not a question of uncertainty about the external world—we are sure of Nate's height. The issue is that the linguistic term “tall” does not refer to a sharp demarcation of objects into two classes—there are *degrees* of tallness. For this reason, *fuzzy set theory is not a method for uncertain reasoning at all*. Rather, fuzzy set theory treats *Tall* as a fuzzy predicate and says that the truth value of *Tall(Nate)* is a number between 0 and 1, rather than being just *true* or *false*. The name “fuzzy set” derives from the interpretation of the predicate as implicitly defining a set of its members—a set that does not have sharp boundaries.

#### FUZZY LOGIC

**Fuzzy logic** is a method for reasoning with logical expressions describing membership in fuzzy sets. For example, the complex sentence  $Tall(Nate) \wedge Heavy(Nate)$  has a fuzzy truth value that is a function of the truth values of its components. The standard rules for evaluating the fuzzy truth,  $T$ , of a complex sentence are

$$\begin{aligned} T(A \wedge B) &= \min(T(A), T(B)) \\ T(A \vee B) &= \max(T(A), T(B)) \\ T(\neg A) &= 1 - T(A) . \end{aligned}$$

Fuzzy logic is therefore a truth-functional system—a fact that causes serious difficulties. For example, suppose that  $T(Tall(Nate)) = 0.6$  and  $T(Heavy(Nate)) = 0.4$ . Then we have  $T(Tall(Nate) \wedge Heavy(Nate)) = 0.4$ , which seems reasonable, but we also get the result  $T(Tall(Nate) \wedge \neg Tall(Nate)) = 0.4$ , which does not. Clearly, the problem arises from the inability of a truth-functional approach to take into account the correlations or anticorrelations among the component propositions.

#### FUZZY CONTROL

**Fuzzy control** is a methodology for constructing control systems in which the mapping between real-valued input and output parameters is represented by fuzzy rules. Fuzzy control has been very successful in commercial products such as automatic transmissions, video

cameras, and electric shavers. Critics (see, e.g., Elkan, 1993) argue that these applications are successful because they have small rule bases, no chaining of inferences, and tunable parameters that can be adjusted to improve the system's performance. The fact that they are implemented with fuzzy operators might be incidental to their success; the key is simply to provide a concise and intuitive way to specify a smoothly interpolated, real-valued function.

There have been attempts to provide an explanation of fuzzy logic in terms of probability theory. One idea is to view assertions such as “Nate is Tall” as discrete observations made concerning a continuous hidden variable, Nate's actual *Height*. The probability model specifies  $P(\text{Observer says Nate is tall} \mid \text{Height})$ , perhaps using a **probit distribution** as described on page 522. A posterior distribution over Nate's height can then be calculated in the usual way, for example, if the model is part of a hybrid Bayesian network. Such an approach is not truth-functional, of course. For example, the conditional distribution

$$P(\text{Observer says Nate is tall and heavy} \mid \text{Height, Weight})$$

allows for interactions between height and weight in the causing of the observation. Thus, someone who is eight feet tall and weighs 190 pounds is very unlikely to be called “tall and heavy,” even though “eight feet” counts as “tall” and “190 pounds” counts as “heavy.”

Fuzzy predicates can also be given a probabilistic interpretation in terms of **random sets**—that is, random variables whose possible values are sets of objects. For example, *Tall* is a random set whose possible values are sets of people. The probability  $P(Tall = S_1)$ , where  $S_1$  is some particular set of people, is the probability that exactly that set would be identified as “tall” by an observer. Then the probability that “Nate is tall” is the sum of the probabilities of all the sets of which Nate is a member.

Both the hybrid Bayesian network approach and the random sets approach appear to capture aspects of fuzziness without introducing degrees of truth. Nonetheless, there remain many open issues concerning the proper representation of linguistic observations and continuous quantities—issues that have been neglected by most outside the fuzzy community.

## 14.8 SUMMARY

This chapter has described **Bayesian networks**, a well-developed representation for uncertain knowledge. Bayesian networks play a role roughly analogous to that of propositional logic for definite knowledge.

- A Bayesian network is a directed acyclic graph whose nodes correspond to random variables; each node has a conditional distribution for the node, given its parents.
- Bayesian networks provide a concise way to represent **conditional independence** relationships in the domain.
- A Bayesian network specifies a full joint distribution; each joint entry is defined as the product of the corresponding entries in the local conditional distributions. A Bayesian network is often exponentially smaller than an explicitly enumerated joint distribution.
- Many conditional distributions can be represented compactly by canonical families of

distributions. **Hybrid Bayesian networks**, which include both discrete and continuous variables, use a variety of canonical distributions.

- Inference in Bayesian networks means computing the probability distribution of a set of query variables, given a set of evidence variables. Exact inference algorithms, such as **variable elimination**, evaluate sums of products of conditional probabilities as efficiently as possible.
- In **polytrees** (singly connected networks), exact inference takes time linear in the size of the network. In the general case, the problem is intractable.
- Stochastic approximation techniques such as **likelihood weighting** and **Markov chain Monte Carlo** can give reasonable estimates of the true posterior probabilities in a network and can cope with much larger networks than can exact algorithms.
- Probability theory can be combined with representational ideas from first-order logic to produce very powerful systems for reasoning under uncertainty. **Relational probability models** (RPMs) include representational restrictions that guarantee a well-defined probability distribution that can be expressed as an equivalent Bayesian network. **Open-universe probability models** handle **existence** and **identity uncertainty**, defining probability distributions over the infinite space of first-order possible worlds.
- Various alternative systems for reasoning under uncertainty have been suggested. Generally speaking, **truth-functional** systems are not well suited for such reasoning.

---

## BIBLIOGRAPHICAL AND HISTORICAL NOTES

The use of networks to represent probabilistic information began early in the 20th century, with the work of Sewall Wright on the probabilistic analysis of genetic inheritance and animal growth factors (Wright, 1921, 1934). I. J. Good (1961), in collaboration with Alan Turing, developed probabilistic representations and Bayesian inference methods that could be regarded as a forerunner of modern Bayesian networks—although the paper is not often cited in this context.<sup>10</sup> The same paper is the original source for the noisy-OR model.

The **influence diagram** representation for decision problems, which incorporated a DAG representation for random variables, was used in decision analysis in the late 1970s (see Chapter 16), but only enumeration was used for evaluation. Judea Pearl developed the message-passing method for carrying out inference in tree networks (Pearl, 1982a) and poly-tree networks (Kim and Pearl, 1983) and explained the importance of causal rather than diagnostic probability models, in contrast to the certainty-factor systems then in vogue.

The first expert system using Bayesian networks was CONVINCER (Kim, 1983). Early applications in medicine included the MUNIN system for diagnosing neuromuscular disorders (Andersen *et al.*, 1989) and the PATHFINDER system for pathology (Heckerman, 1991). The CPCS system (Pradhan *et al.*, 1994) is a Bayesian network for internal medicine consisting

---

<sup>10</sup> I. J. Good was chief statistician for Turing's code-breaking team in World War II. In *2001: A Space Odyssey* (Clarke, 1968a), Good and Minsky are credited with making the breakthrough that led to the development of the HAL 9000 computer.

of 448 nodes, 906 links and 8,254 conditional probability values. (The front cover shows a portion of the network.)

Applications in engineering include the Electric Power Research Institute's work on monitoring power generators (Morjaria *et al.*, 1995), NASA's work on displaying time-critical information at Mission Control in Houston (Horvitz and Barry, 1995), and the general field of **network tomography**, which aims to infer unobserved local properties of nodes and links in the Internet from observations of end-to-end message performance (Castro *et al.*, 2004). Perhaps the most widely used Bayesian network systems have been the diagnosis-and-repair modules (e.g., the Printer Wizard) in Microsoft Windows (Breese and Heckerman, 1996) and the Office Assistant in Microsoft Office (Horvitz *et al.*, 1998). Another important application area is biology: Bayesian networks have been used for identifying human genes by reference to mouse genes (Zhang *et al.*, 2003), inferring cellular networks Friedman (2004), and many other tasks in bioinformatics. We could go on, but instead we'll refer you to Pourret *et al.* (2008), a 400-page guide to applications of Bayesian networks.

Ross Shachter (1986), working in the influence diagram community, developed the first complete algorithm for general Bayesian networks. His method was based on goal-directed reduction of the network using posterior-preserving transformations. Pearl (1986) developed a clustering algorithm for exact inference in general Bayesian networks, utilizing a conversion to a directed polytree of clusters in which message passing was used to achieve consistency over variables shared between clusters. A similar approach, developed by the statisticians David Spiegelhalter and Steffen Lauritzen (Lauritzen and Spiegelhalter, 1988), is based on conversion to an undirected form of graphical model called a **Markov network**. This approach is implemented in the HUGIN system, an efficient and widely used tool for uncertain reasoning (Andersen *et al.*, 1989). Boutilier *et al.* (1996) show how to exploit context-specific independence in clustering algorithms.

The basic idea of variable elimination—that repeated computations within the overall sum-of-products expression can be avoided by caching—appeared in the symbolic probabilistic inference (SPI) algorithm (Shachter *et al.*, 1990). The elimination algorithm we describe is closest to that developed by Zhang and Poole (1994). Criteria for pruning irrelevant variables were developed by Geiger *et al.* (1990) and by Lauritzen *et al.* (1990); the criterion we give is a simple special case of these. Dechter (1999) shows how the variable elimination idea is essentially identical to **nonserial dynamic programming** (Bertele and Brioschi, 1972), an algorithmic approach that can be applied to solve a range of inference problems in Bayesian networks—for example, finding the **most likely explanation** for a set of observations. This connects Bayesian network algorithms to related methods for solving CSPs and gives a direct measure of the complexity of exact inference in terms of the tree width of the network. Wexler and Meek (2009) describe a method of preventing exponential growth in the size of factors computed in variable elimination; their algorithm breaks down large factors into products of smaller factors and simultaneously computes an error bound for the resulting approximation.

The inclusion of continuous random variables in Bayesian networks was considered by Pearl (1988) and Shachter and Kenley (1989); these papers discussed networks containing only continuous variables with linear Gaussian distributions. The inclusion of discrete variables has been investigated by Lauritzen and Wermuth (1989) and implemented in the

MARKOV NETWORK

NONSERIAL DYNAMIC  
PROGRAMMING

cHUGIN system (Olesen, 1993). Further analysis of linear Gaussian models, with connections to many other models used in statistics, appears in Roweis and Ghahramani (1999). The probit distribution is usually attributed to Gaddum (1933) and Bliss (1934), although it had been discovered several times in the 19th century. Bliss's work was expanded considerably by Finney (1947). The probit has been used widely for modeling discrete choice phenomena and can be extended to handle more than two choices (Daganzo, 1979). The logit model was introduced by Berkson (1944); initially much derided, it eventually became more popular than the probit model. Bishop (1995) gives a simple justification for its use.

Cooper (1990) showed that the general problem of inference in unconstrained Bayesian networks is NP-hard, and Paul Dagum and Mike Luby (1993) showed the corresponding approximation problem to be NP-hard. Space complexity is also a serious problem in both clustering and variable elimination methods. The method of **cutset conditioning**, which was developed for CSPs in Chapter 6, avoids the construction of exponentially large tables. In a Bayesian network, a cutset is a set of nodes that, when instantiated, reduces the remaining nodes to a polytree that can be solved in linear time and space. The query is answered by summing over all the instantiations of the cutset, so the overall space requirement is still linear (Pearl, 1988). Darwiche (2001) describes a recursive conditioning algorithm that allows a complete range of space/time tradeoffs.

The development of fast approximation algorithms for Bayesian network inference is a very active area, with contributions from statistics, computer science, and physics. The rejection sampling method is a general technique that is long known to statisticians; it was first applied to Bayesian networks by Max Henrion (1988), who called it **logic sampling**. Likelihood weighting, which was developed by Fung and Chang (1989) and Shachter and Peot (1989), is an example of the well-known statistical method of **importance sampling**. Cheng and Druzdzel (2000) describe an adaptive version of likelihood weighting that works well even when the evidence has very low prior likelihood.

Markov chain Monte Carlo (MCMC) algorithms began with the Metropolis algorithm, due to Metropolis *et al.* (1953), which was also the source of the simulated annealing algorithm described in Chapter 4. The Gibbs sampler was devised by Geman and Geman (1984) for inference in undirected Markov networks. The application of MCMC to Bayesian networks is due to Pearl (1987). The papers collected by Gilks *et al.* (1996) cover a wide variety of applications of MCMC, several of which were developed in the well-known BUGS package (Gilks *et al.*, 1994).

There are two very important families of approximation methods that we did not cover in the chapter. The first is the family of **variational approximation** methods, which can be used to simplify complex calculations of all kinds. The basic idea is to propose a reduced version of the original problem that is simple to work with, but that resembles the original problem as closely as possible. The reduced problem is described by some **variational parameters**  $\lambda$  that are adjusted to minimize a distance function  $D$  between the original and the reduced problem, often by solving the system of equations  $\partial D / \partial \lambda = 0$ . In many cases, strict upper and lower bounds can be obtained. Variational methods have long been used in statistics (Rustagi, 1976). In statistical physics, the **mean-field** method is a particular variational approximation in which the individual variables making up the model are assumed

VARIATIONAL  
APPROXIMATION

VARIATIONAL  
PARAMETER

MEAN FIELD

to be completely independent. This idea was applied to solve large undirected Markov networks (Peterson and Anderson, 1987; Parisi, 1988). Saul *et al.* (1996) developed the mathematical foundations for applying variational methods to Bayesian networks and obtained accurate lower-bound approximations for sigmoid networks with the use of mean-field methods. Jaakkola and Jordan (1996) extended the methodology to obtain both lower and upper bounds. Since these early papers, variational methods have been applied to many specific families of models. The remarkable paper by Wainwright and Jordan (2008) provides a unifying theoretical analysis of the literature on variational methods.

A second important family of approximation algorithms is based on Pearl's polytree message-passing algorithm (1982a). This algorithm can be applied to general networks, as suggested by Pearl (1988). The results might be incorrect, or the algorithm might fail to terminate, but in many cases, the values obtained are close to the true values. Little attention was paid to this so-called **belief propagation** (or BP) approach until McEliece *et al.* (1998) observed that message passing in a multiply connected Bayesian network was exactly the computation performed by the **turbo decoding** algorithm (Berrou *et al.*, 1993), which provided a major breakthrough in the design of efficient error-correcting codes. The implication is that BP is both fast and accurate on the very large and very highly connected networks used for decoding and might therefore be useful more generally. Murphy *et al.* (1999) presented a promising empirical study of BP's performance, and Weiss and Freeman (2001) established strong convergence results for BP on linear Gaussian networks. Weiss (2000b) shows how an approximation called loopy belief propagation works, and when the approximation is correct. Yedidia *et al.* (2005) made further connections between loopy propagation and ideas from statistical physics.

The connection between probability and first-order languages was first studied by Carnap (1950). Gaifman (1964) and Scott and Krauss (1966) defined a language in which probabilities could be associated with first-order sentences and for which models were probability measures on possible worlds. Within AI, this idea was developed for propositional logic by Nilsson (1986) and for first-order logic by Halpern (1990). The first extensive investigation of knowledge representation issues in such languages was carried out by Bacchus (1990). The basic idea is that each sentence in the knowledge base expressed a *constraint* on the distribution over possible worlds; one sentence entails another if it expresses a stronger constraint. For example, the sentence  $\forall x \ P(Hungry(x)) > 0.2$  rules out distributions in which any object is hungry with probability less than 0.2; thus, it entails the sentence  $\forall x \ P(Hungry(x)) > 0.1$ . It turns out that writing a *consistent* set of sentences in these languages is quite difficult and constructing a unique probability model nearly impossible unless one adopts the representation approach of Bayesian networks by writing suitable sentences about conditional probabilities.

Beginning in the early 1990s, researchers working on complex applications noticed the expressive limitations of Bayesian networks and developed various languages for writing "templates" with logical variables, from which large networks could be constructed automatically for each problem instance (Breese, 1992; Wellman *et al.*, 1992). The most important such language was BUGS (Bayesian inference Using Gibbs Sampling) (Gilks *et al.*, 1994), which combined Bayesian networks with the **indexed random variable** notation common in

BELIEF  
PROPAGATION

TURBO DECODING

INDEXED RANDOM  
VARIABLE

statistics. (In BUGS, an indexed random variable looks like  $X[i]$ , where  $i$  has a defined integer range.) These languages inherited the key property of Bayesian networks: every well-formed knowledge base defines a unique, consistent probability model. Languages with well-defined semantics based on unique names and domain closure drew on the representational capabilities of logic programming (Poole, 1993; Sato and Kameya, 1997; Kersting *et al.*, 2000) and semantic networks (Koller and Pfeffer, 1998; Pfeffer, 2000). Pfeffer (2007) went on to develop IBAL, which represents first-order probability models as probabilistic programs in a programming language extended with a randomization primitive. Another important thread was the combination of relational and first-order notations with (undirected) Markov networks (Taskar *et al.*, 2002; Domingos and Richardson, 2004), where the emphasis has been less on knowledge representation and more on learning from large data sets.

Initially, inference in these models was performed by generating an equivalent Bayesian network. Pfeffer *et al.* (1999) introduced a variable elimination algorithm that cached each computed factor for reuse by later computations involving the same relations but different objects, thereby realizing some of the computational gains of lifting. The first truly lifted inference algorithm was a lifted form of variable elimination described by Poole (2003) and subsequently improved by de Salvo Braz *et al.* (2007). Further advances, including cases where certain aggregate probabilities can be computed in closed form, are described by Milch *et al.* (2008) and Kisynski and Poole (2009). Pasula and Russell (2001) studied the application of MCMC to avoid building the complete equivalent Bayes net in cases of relational and identity uncertainty. Getoor and Taskar (2007) collect many important papers on first-order probability models and their use in machine learning.

#### RECORD LINKAGE

Probabilistic reasoning about identity uncertainty has two distinct origins. In statistics, the problem of **record linkage** arises when data records do not contain standard unique identifiers—for example, various citations of this book might name its first author “Stuart Russell” or “S. J. Russell” or even “Stewart Russle,” and other authors may use the some of the same names. Literally hundreds of companies exist solely to solve record linkage problems in financial, medical, census, and other data. Probabilistic analysis goes back to work by Dunn (1946); the Fellegi–Sunter model (1969), which is essentially naive Bayes applied to matching, still dominates current practice. The second origin for work on identity uncertainty is multitarget tracking (Sittler, 1964), which we cover in Chapter 15. For most of its history, work in symbolic AI assumed erroneously that sensors could supply sentences with unique identifiers for objects. The issue was studied in the context of language understanding by Charniak and Goldman (1992) and in the context of surveillance by (Huang and Russell, 1998) and Pasula *et al.* (1999). Pasula *et al.* (2003) developed a complex generative model for authors, papers, and citation strings, involving both relational and identity uncertainty, and demonstrated high accuracy for citation information extraction. The first formally defined language for open-universe probability models was BLOG (Milch *et al.*, 2005), which came with a complete (albeit slow) MCMC inference algorithm for all well-defined models. (The program code faintly visible on the front cover of this book is part of a BLOG model for detecting nuclear explosions from seismic signals as part of the UN Comprehensive Test Ban Treaty verification regime.) Laskey (2008) describes another open-universe modeling language called **multi-entity Bayesian networks**.



As explained in Chapter 13, early probabilistic systems fell out of favor in the early 1970s, leaving a partial vacuum to be filled by alternative methods. Certainty factors were invented for use in the medical expert system MYCIN (Shortliffe, 1976), which was intended both as an engineering solution and as a model of human judgment under uncertainty. The collection *Rule-Based Expert Systems* (Buchanan and Shortliffe, 1984) provides a complete overview of MYCIN and its descendants (see also Stefik, 1995). David Heckerman (1986) showed that a slightly modified version of certainty factor calculations gives correct probabilistic results in some cases, but results in serious overcounting of evidence in other cases. The PROSPECTOR expert system (Duda *et al.*, 1979) used a rule-based approach in which the rules were justified by a (seldom tenable) global independence assumption.

Dempster–Shafer theory originates with a paper by Arthur Dempster (1968) proposing a generalization of probability to interval values and a combination rule for using them. Later work by Glenn Shafer (1976) led to the Dempster–Shafer theory’s being viewed as a competing approach to probability. Pearl (1988) and Ruspini *et al.* (1992) analyze the relationship between the Dempster–Shafer theory and standard probability theory.

Fuzzy sets were developed by Lotfi Zadeh (1965) in response to the perceived difficulty of providing exact inputs to intelligent systems. The text by Zimmermann (2001) provides a thorough introduction to fuzzy set theory; papers on fuzzy applications are collected in Zimmermann (1999). As we mentioned in the text, fuzzy logic has often been perceived incorrectly as a direct competitor to probability theory, whereas in fact it addresses a different set of issues. **Possibility theory** (Zadeh, 1978) was introduced to handle uncertainty in fuzzy systems and has much in common with probability. Dubois and Prade (1994) survey the connections between possibility theory and probability theory.

The resurgence of probability depended mainly on Pearl’s development of Bayesian networks as a method for representing and using conditional independence information. This resurgence did not come without a fight; Peter Cheeseman’s (1985) pugnacious “In Defense of Probability” and his later article “An Inquiry into Computer Understanding” (Cheeseman, 1988, with commentaries) give something of the flavor of the debate. Eugene Charniak helped present the ideas to AI researchers with a popular article, “Bayesian networks without tears”<sup>11</sup> (1991), and book (1993). The book by Dean and Wellman (1991) also helped introduce Bayesian networks to AI researchers. One of the principal philosophical objections of the logicians was that the numerical calculations that probability theory was thought to require were not apparent to introspection and presumed an unrealistic level of precision in our uncertain knowledge. The development of **qualitative probabilistic networks** (Wellman, 1990a) provided a purely qualitative abstraction of Bayesian networks, using the notion of positive and negative influences between variables. Wellman shows that in many cases such information is sufficient for optimal decision making without the need for the precise specification of probability values. Goldszmidt and Pearl (1996) take a similar approach. Work by Adnan Darwiche and Matt Ginsberg (1992) extracts the basic properties of conditioning and evidence combination from probability theory and shows that they can also be applied in logical and default reasoning. Often, programs speak louder than words, and the ready avail-

<sup>11</sup> The title of the original version of the article was “Pearl for swine.”

ability of high-quality software such as the Bayes Net toolkit (Murphy, 2001) accelerated the adoption of the technology.

The most important single publication in the growth of Bayesian networks was undoubtedly the text *Probabilistic Reasoning in Intelligent Systems* (Pearl, 1988). Several excellent texts (Lauritzen, 1996; Jensen, 2001; Korb and Nicholson, 2003; Jensen, 2007; Darwiche, 2009; Koller and Friedman, 2009) provide thorough treatments of the topics we have covered in this chapter. New research on probabilistic reasoning appears both in mainstream AI journals, such as *Artificial Intelligence* and the *Journal of AI Research*, and in more specialized journals, such as the *International Journal of Approximate Reasoning*. Many papers on graphical models, which include Bayesian networks, appear in statistical journals. The proceedings of the conferences on Uncertainty in Artificial Intelligence (UAI), Neural Information Processing Systems (NIPS), and Artificial Intelligence and Statistics (AISTATS) are excellent sources for current research.

---

## EXERCISES

**14.1** We have a bag of three biased coins  $a$ ,  $b$ , and  $c$  with probabilities of coming up heads of 20%, 60%, and 80%, respectively. One coin is drawn randomly from the bag (with equal likelihood of drawing each of the three coins), and then the coin is flipped three times to generate the outcomes  $X_1$ ,  $X_2$ , and  $X_3$ .

- a. Draw the Bayesian network corresponding to this setup and define the necessary CPTs.
- b. Calculate which coin was most likely to have been drawn from the bag if the observed flips come out heads twice and tails once.

**14.2** Equation (14.1) on page 513 defines the joint distribution represented by a Bayesian network in terms of the parameters  $\theta(X_i | \text{Parents}(X_i))$ . This exercise asks you to derive the equivalence between the parameters and the conditional probabilities  $\mathbf{P}(X_i | \text{Parents}(X_i))$  from this definition.

- a. Consider a simple network  $X \rightarrow Y \rightarrow Z$  with three Boolean variables. Use Equations (13.3) and (13.6) (pages 485 and 492) to express the conditional probability  $P(z | y)$  as the ratio of two sums, each over entries in the joint distribution  $\mathbf{P}(X, Y, Z)$ .
- b. Now use Equation (14.1) to write this expression in terms of the network parameters  $\theta(X)$ ,  $\theta(Y | X)$ , and  $\theta(Z | Y)$ .
- c. Next, expand out the summations in your expression from part (b), writing out explicitly the terms for the true and false values of each summed variable. Assuming that all network parameters satisfy the constraint  $\sum_{x_i} \theta(x_i | \text{parents}(X_i)) = 1$ , show that the resulting expression reduces to  $\theta(x | y)$ .
- d. Generalize this derivation to show that  $\theta(X_i | \text{Parents}(X_i)) = \mathbf{P}(X_i | \text{Parents}(X_i))$  for any Bayesian network.

## ARC REVERSAL

**14.3** The operation of **arc reversal** in a Bayesian network allows us to change the direction of an arc  $X \rightarrow Y$  while preserving the joint probability distribution that the network represents (Shachter, 1986). Arc reversal may require introducing new arcs: all the parents of  $X$  also become parents of  $Y$ , and all parents of  $Y$  also become parents of  $X$ .

- Assume that  $X$  and  $Y$  start with  $m$  and  $n$  parents, respectively, and that all variables have  $k$  values. By calculating the change in size for the CPTs of  $X$  and  $Y$ , show that the total number of parameters in the network cannot decrease during arc reversal. (*Hint: the parents of  $X$  and  $Y$  need not be disjoint.*)
- Under what circumstances can the total number remain constant?
- Let the parents of  $X$  be  $\mathbf{U} \cup \mathbf{V}$  and the parents of  $Y$  be  $\mathbf{V} \cup \mathbf{W}$ , where  $\mathbf{U}$  and  $\mathbf{W}$  are disjoint. The formulas for the new CPTs after arc reversal are as follows:

$$\mathbf{P}(Y | \mathbf{U}, \mathbf{V}, \mathbf{W}) = \sum_x \mathbf{P}(Y | \mathbf{V}, \mathbf{W}, x) \mathbf{P}(x | \mathbf{U}, \mathbf{V})$$

$$\mathbf{P}(X | \mathbf{U}, \mathbf{V}, \mathbf{W}, Y) = \mathbf{P}(Y | X, \mathbf{V}, \mathbf{W}) \mathbf{P}(X | \mathbf{U}, \mathbf{V}) / \mathbf{P}(Y | \mathbf{U}, \mathbf{V}, \mathbf{W}) .$$

Prove that the new network expresses the same joint distribution over all variables as the original network.

**14.4** Consider the Bayesian network in Figure 14.2.

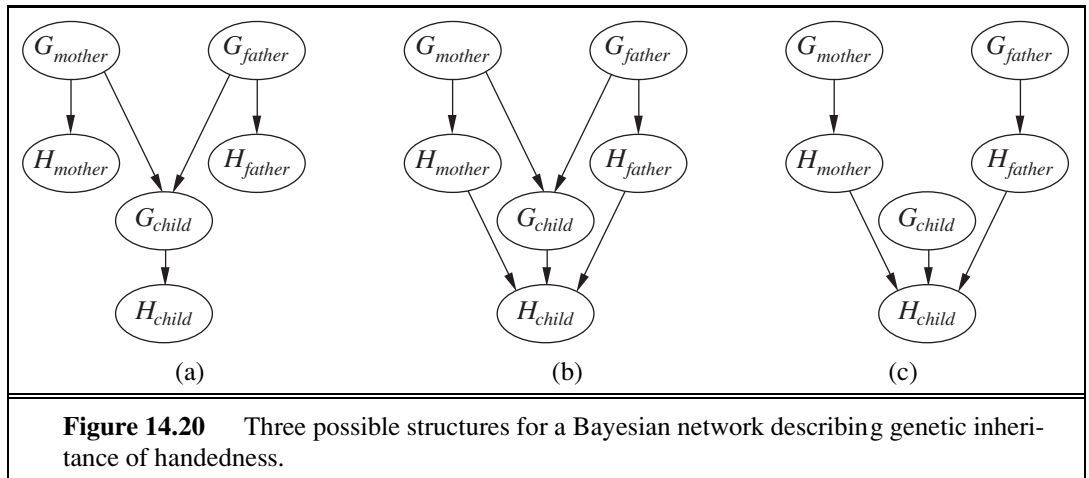
- If no evidence is observed, are *Burglary* and *Earthquake* independent? Prove this from the numerical semantics and from the topological semantics.
- If we observe  $Alarm = true$ , are *Burglary* and *Earthquake* independent? Justify your answer by calculating whether the probabilities involved satisfy the definition of conditional independence.

**14.5** Suppose that in a Bayesian network containing an unobserved variable  $Y$ , all the variables in the Markov blanket  $MB(Y)$  have been observed.

- Prove that removing the node  $Y$  from the network will not affect the posterior distribution for any other unobserved variable in the network.
- Discuss whether we can remove  $Y$  if we are planning to use (i) rejection sampling and (ii) likelihood weighting.

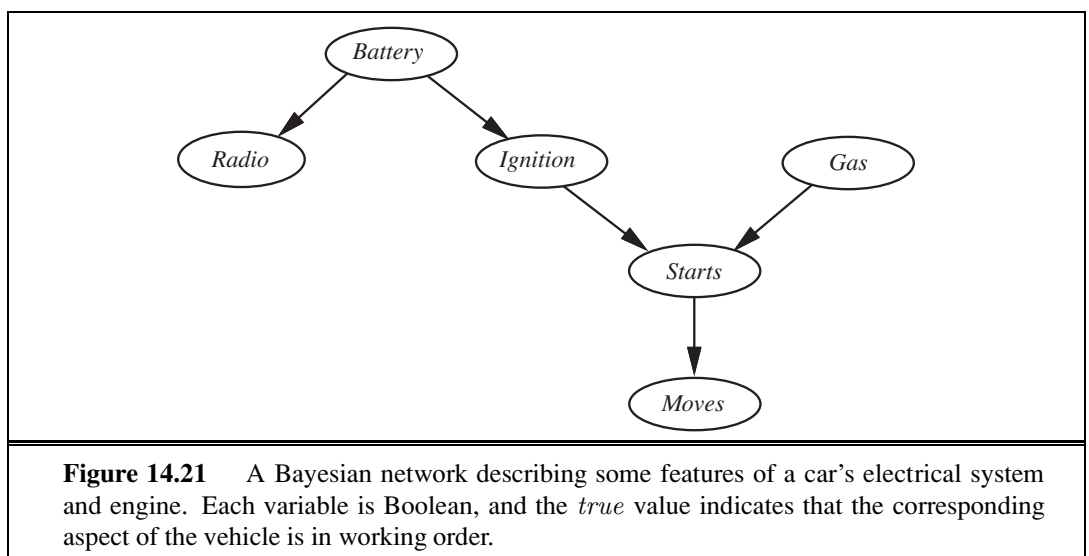
**14.6** Let  $H_x$  be a random variable denoting the handedness of an individual  $x$ , with possible values  $l$  or  $r$ . A common hypothesis is that left- or right-handedness is inherited by a simple mechanism; that is, perhaps there is a gene  $G_x$ , also with values  $l$  or  $r$ , and perhaps actual handedness turns out mostly the same (with some probability  $s$ ) as the gene an individual possesses. Furthermore, perhaps the gene itself is equally likely to be inherited from either of an individual's parents, with a small nonzero probability  $m$  of a random mutation flipping the handedness.

- Which of the three networks in Figure 14.20 claim that  $\mathbf{P}(G_{father}, G_{mother}, G_{child}) = \mathbf{P}(G_{father})\mathbf{P}(G_{mother})\mathbf{P}(G_{child})$ ?
- Which of the three networks make independence claims that are consistent with the hypothesis about the inheritance of handedness?



- c. Which of the three networks is the best description of the hypothesis?
- d. Write down the CPT for the  $G_{child}$  node in network (a), in terms of  $s$  and  $m$ .
- e. Suppose that  $P(G_{father} = l) = P(G_{mother} = l) = q$ . In network (a), derive an expression for  $P(G_{child} = l)$  in terms of  $m$  and  $q$  only, by conditioning on its parent nodes.
- f. Under conditions of genetic equilibrium, we expect the distribution of genes to be the same across generations. Use this to calculate the value of  $q$ , and, given what you know about handedness in humans, explain why the hypothesis described at the beginning of this question must be wrong.

**14.7** The **Markov blanket** of a variable is defined on page 517. Prove that a variable is independent of all other variables in the network, given its Markov blanket and derive Equation (14.12) (page 538).



**14.8** Consider the network for car diagnosis shown in Figure 14.21.

- a. Extend the network with the Boolean variables *IcyWeather* and *StarterMotor*.
- b. Give reasonable conditional probability tables for all the nodes.
- c. How many independent values are contained in the joint probability distribution for eight Boolean nodes, assuming that no conditional independence relations are known to hold among them?
- d. How many independent probability values do your network tables contain?
- e. The conditional distribution for *Starts* could be described as a **noisy-AND** distribution. Define this family in general and relate it to the noisy-OR distribution.

**14.9** Consider the family of linear Gaussian networks, as defined on page 520.

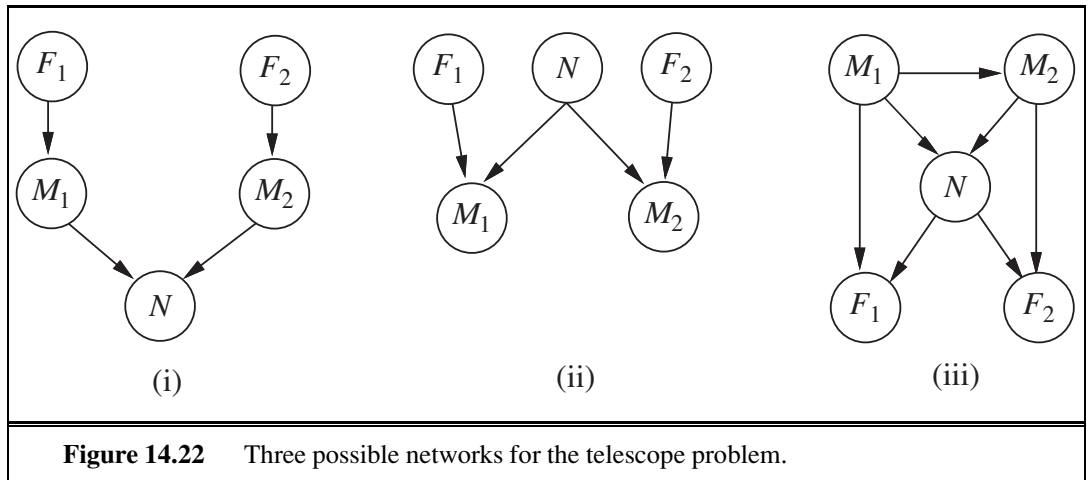
- a. In a two-variable network, let  $X_1$  be the parent of  $X_2$ , let  $X_1$  have a Gaussian prior, and let  $\mathbf{P}(X_2 | X_1)$  be a linear Gaussian distribution. Show that the joint distribution  $P(X_1, X_2)$  is a multivariate Gaussian, and calculate its covariance matrix.
- b. Prove by induction that the joint distribution for a general linear Gaussian network on  $X_1, \dots, X_n$  is also a multivariate Gaussian.

**14.10** The probit distribution defined on page 522 describes the probability distribution for a Boolean child, given a single continuous parent.

- a. How might the definition be extended to cover multiple continuous parents?
- b. How might it be extended to handle a *multivalued* child variable? Consider both cases where the child's values are ordered (as in selecting a gear while driving, depending on speed, slope, desired acceleration, etc.) and cases where they are unordered (as in selecting bus, train, or car to get to work). (*Hint*: Consider ways to divide the possible values into two sets, to mimic a Boolean variable.)

**14.11** In your local nuclear power station, there is an alarm that senses when a temperature gauge exceeds a given threshold. The gauge measures the temperature of the core. Consider the Boolean variables  $A$  (alarm sounds),  $F_A$  (alarm is faulty), and  $F_G$  (gauge is faulty) and the multivalued nodes  $G$  (gauge reading) and  $T$  (actual core temperature).

- a. Draw a Bayesian network for this domain, given that the gauge is more likely to fail when the core temperature gets too high.
- b. Is your network a polytree? Why or why not?
- c. Suppose there are just two possible actual and measured temperatures, normal and high; the probability that the gauge gives the correct temperature is  $x$  when it is working, but  $y$  when it is faulty. Give the conditional probability table associated with  $G$ .
- d. Suppose the alarm works correctly unless it is faulty, in which case it never sounds. Give the conditional probability table associated with  $A$ .
- e. Suppose the alarm and gauge are working and the alarm sounds. Calculate an expression for the probability that the temperature of the core is too high, in terms of the various conditional probabilities in the network.



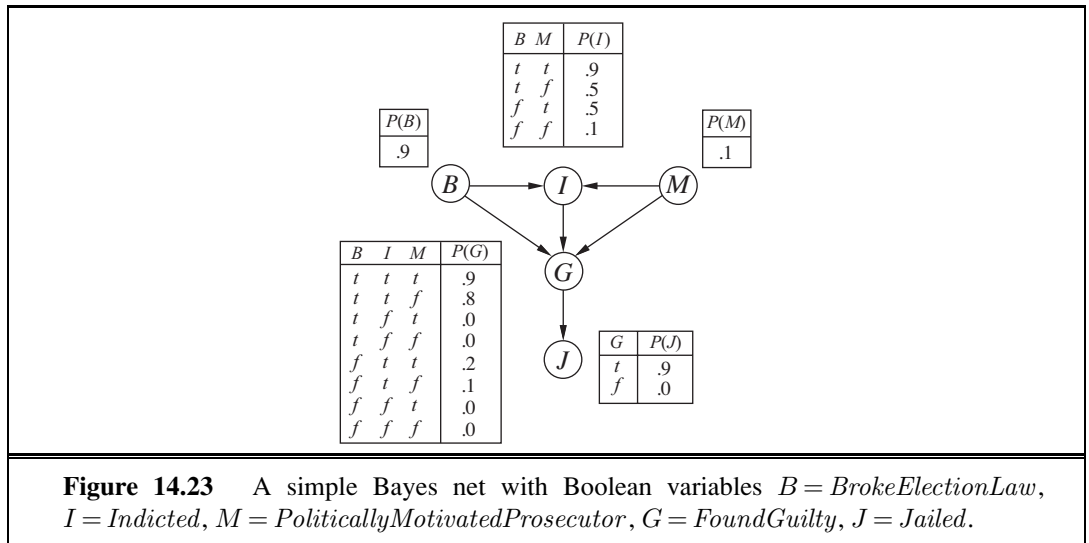
**14.12** Two astronomers in different parts of the world make measurements  $M_1$  and  $M_2$  of the number of stars  $N$  in some small region of the sky, using their telescopes. Normally, there is a small possibility  $e$  of error by up to one star in each direction. Each telescope can also (with a much smaller probability  $f$ ) be badly out of focus (events  $F_1$  and  $F_2$ ), in which case the scientist will undercount by three or more stars (or if  $N$  is less than 3, fail to detect any stars at all). Consider the three networks shown in Figure 14.22.

- Which of these Bayesian networks are correct (but not necessarily efficient) representations of the preceding information?
- Which is the best network? Explain.
- Write out a conditional distribution for  $\mathbf{P}(M_1 \mid N)$ , for the case where  $N \in \{1, 2, 3\}$  and  $M_1 \in \{0, 1, 2, 3, 4\}$ . Each entry in the conditional distribution should be expressed as a function of the parameters  $e$  and/or  $f$ .
- Suppose  $M_1 = 1$  and  $M_2 = 3$ . What are the *possible* numbers of stars if you assume no prior constraint on the values of  $N$ ?
- What is the *most likely* number of stars, given these observations? Explain how to compute this, or if it is not possible to compute, explain what additional information is needed and how it would affect the result.

**14.13** Consider the network shown in Figure 14.22(ii), and assume that the two telescopes work identically.  $N \in \{1, 2, 3\}$  and  $M_1, M_2 \in \{0, 1, 2, 3, 4\}$ , with the symbolic CPTs as described in Exercise 14.12. Using the enumeration algorithm (Figure 14.9 on page 525), calculate the probability distribution  $\mathbf{P}(N \mid M_1 = 2, M_2 = 2)$ .

**14.14** Consider the Bayes net shown in Figure 14.23.

- Which of the following are asserted by the network *structure*?
  - $\mathbf{P}(B, I, M) = \mathbf{P}(B)\mathbf{P}(I)\mathbf{P}(M)$ .
  - $\mathbf{P}(J \mid G) = \mathbf{P}(J \mid G, I)$ .
  - $\mathbf{P}(M \mid G, B, I) = \mathbf{P}(M \mid G, B, I, J)$ .



- Calculate the value of  $P(b, i, \neg m, g, j)$ .
- Calculate the probability that someone goes to jail given that they broke the law, have been indicted, and face a politically motivated prosecutor.
- A **context-specific independence** (see page 542) allows a variable to be independent of some of its parents given certain values of others. In addition to the usual conditional independences given by the graph structure, what context-specific independences exist in the Bayes net in Figure 14.23?
- Suppose we want to add the variable  $P = \text{PresidentialPardon}$  to the network; draw the new network and briefly explain any links you add.

**14.15** Consider the variable elimination algorithm in Figure 14.11 (page 528).

- Section 14.4 applies variable elimination to the query

$$\mathbf{P}(\text{Burglary} \mid \text{JohnCalls} = \text{true}, \text{MaryCalls} = \text{true}) .$$

Perform the calculations indicated and check that the answer is correct.

- Count the number of arithmetic operations performed, and compare it with the number performed by the enumeration algorithm.
- Suppose a network has the form of a *chain*: a sequence of Boolean variables  $X_1, \dots, X_n$  where  $\text{Parents}(X_i) = \{X_{i-1}\}$  for  $i = 2, \dots, n$ . What is the complexity of computing  $\mathbf{P}(X_1 \mid X_n = \text{true})$  using enumeration? Using variable elimination?
- Prove that the complexity of running variable elimination on a polytree network is linear in the size of the tree for any variable ordering consistent with the network structure.

**14.16** Investigate the complexity of exact inference in general Bayesian networks:

- Prove that any 3-SAT problem can be reduced to exact inference in a Bayesian network constructed to represent the particular problem and hence that exact inference is NP-

hard. (*Hint*: Consider a network with one variable for each proposition symbol, one for each clause, and one for the conjunction of clauses.)

- b. The problem of counting the number of satisfying assignments for a 3-SAT problem is #P-complete. Show that exact inference is at least as hard as this.

**14.17** Consider the problem of generating a random sample from a specified distribution on a single variable. Assume you have a random number generator that returns a random number uniformly distributed between 0 and 1.

- a. Let  $X$  be a discrete variable with  $P(X = x_i) = p_i$  for  $i \in \{1, \dots, k\}$ . The **cumulative distribution** of  $X$  gives the probability that  $X \in \{x_1, \dots, x_j\}$  for each possible  $j$ . (See also Appendix A.) Explain how to calculate the cumulative distribution in  $O(k)$  time and how to generate a single sample of  $X$  from it. Can the latter be done in less than  $O(k)$  time?
- b. Now suppose we want to generate  $N$  samples of  $X$ , where  $N \gg k$ . Explain how to do this with an expected run time per sample that is *constant* (i.e., independent of  $k$ ).
- c. Now consider a continuous-valued variable with a parameterized distribution (e.g., Gaussian). How can samples be generated from such a distribution?
- d. Suppose you want to query a continuous-valued variable and you are using a sampling algorithm such as LIKELIHOODWEIGHTING to do the inference. How would you have to modify the query-answering process?

**14.18** Consider the query  $\mathbf{P}(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$  in Figure 14.12(a) (page 529) and how Gibbs sampling can answer it.

- a. How many states does the Markov chain have?
- b. Calculate the **transition matrix**  $\mathbf{Q}$  containing  $q(\mathbf{y} \rightarrow \mathbf{y}')$  for all  $\mathbf{y}, \mathbf{y}'$ .
- c. What does  $\mathbf{Q}^2$ , the square of the transition matrix, represent?
- d. What about  $\mathbf{Q}^n$  as  $n \rightarrow \infty$ ?
- e. Explain how to do probabilistic inference in Bayesian networks, assuming that  $\mathbf{Q}^n$  is available. Is this a practical way to do inference?

**14.19** This exercise explores the stationary distribution for Gibbs sampling methods.

- a. The convex composition  $[\alpha, q_1; 1 - \alpha, q_2]$  of  $q_1$  and  $q_2$  is a transition probability distribution that first chooses one of  $q_1$  and  $q_2$  with probabilities  $\alpha$  and  $1 - \alpha$ , respectively, and then applies whichever is chosen. Prove that if  $q_1$  and  $q_2$  are in detailed balance with  $\pi$ , then their convex composition is also in detailed balance with  $\pi$ . (*Note*: this result justifies a variant of GIBBS-ASK in which variables are chosen at random rather than sampled in a fixed sequence.)
- b. Prove that if each of  $q_1$  and  $q_2$  has  $\pi$  as its stationary distribution, then the sequential composition  $q = q_1 \circ q_2$  also has  $\pi$  as its stationary distribution.

**14.20** The **Metropolis–Hastings** algorithm is a member of the MCMC family; as such, it is designed to generate samples  $\mathbf{x}$  (eventually) according to target probabilities  $\pi(\mathbf{x})$ . (Typically

CUMULATIVE  
DISTRIBUTION

METROPOLIS-  
HASTINGS



PROPOSAL  
DISTRIBUTIONACCEPTANCE  
PROBABILITY

we are interested in sampling from  $\pi(\mathbf{x}) = P(\mathbf{x} | \mathbf{e})$ .) Like simulated annealing, Metropolis–Hastings operates in two stages. First, it samples a new state  $\mathbf{x}'$  from a **proposal distribution**  $q(\mathbf{x}' | \mathbf{x})$ , given the current state  $\mathbf{x}$ . Then, it probabilistically accepts or rejects  $\mathbf{x}'$  according to the **acceptance probability**

$$\alpha(\mathbf{x}' | \mathbf{x}) = \min \left( 1, \frac{\pi(\mathbf{x}')q(\mathbf{x} | \mathbf{x}')}{\pi(\mathbf{x})q(\mathbf{x}' | \mathbf{x})} \right).$$

If the proposal is rejected, the state remains at  $\mathbf{x}$ .

- a. Consider an ordinary Gibbs sampling step for a specific variable  $X_i$ . Show that this step, considered as a proposal, is guaranteed to be accepted by Metropolis–Hastings. (Hence, Gibbs sampling is a special case of Metropolis–Hastings.)
- b. Show that the two-step process above, viewed as a transition probability distribution, is in detailed balance with  $\pi$ .



**14.21** Three soccer teams  $A$ ,  $B$ , and  $C$ , play each other once. Each match is between two teams, and can be won, drawn, or lost. Each team has a fixed, unknown degree of quality—an integer ranging from 0 to 3—and the outcome of a match depends probabilistically on the difference in quality between the two teams.

- a. Construct a relational probability model to describe this domain, and suggest numerical values for all the necessary probability distributions.
- b. Construct the equivalent Bayesian network for the three matches.
- c. Suppose that in the first two matches  $A$  beats  $B$  and draws with  $C$ . Using an exact inference algorithm of your choice, compute the posterior distribution for the outcome of the third match.
- d. Suppose there are  $n$  teams in the league and we have the results for all but the last match. How does the complexity of predicting the last game vary with  $n$ ?
- e. Investigate the application of MCMC to this problem. How quickly does it converge in practice and how well does it scale?

# 15 PROBABILISTIC REASONING OVER TIME

*In which we try to interpret the present, understand the past, and perhaps predict the future, even when very little is crystal clear.*

Agents in partially observable environments must be able to keep track of the current state, to the extent that their sensors allow. In Section 4.4 we showed a methodology for doing that: an agent maintains a **belief state** that represents which states of the world are currently possible. From the belief state and a **transition model**, the agent can predict how the world might evolve in the next time step. From the percepts observed and a **sensor model**, the agent can update the belief state. This is a pervasive idea: in Chapter 4 belief states were represented by explicitly enumerated sets of states, whereas in Chapters 7 and 11 they were represented by logical formulas. Those approaches defined belief states in terms of which world states were *possible*, but could say nothing about which states were *likely* or *unlikely*. In this chapter, we use probability theory to quantify the degree of belief in elements of the belief state.

As we show in Section 15.1, time itself is handled in the same way as in Chapter 7: a changing world is modeled using a variable for each aspect of the world state *at each point in time*. The transition and sensor models may be uncertain: the transition model describes the probability distribution of the variables at time  $t$ , given the state of the world at past times, while the sensor model describes the probability of each percept at time  $t$ , given the current state of the world. Section 15.2 defines the basic inference tasks and describes the general structure of inference algorithms for temporal models. Then we describe three specific kinds of models: **hidden Markov models**, **Kalman filters**, and **dynamic Bayesian networks** (which include hidden Markov models and Kalman filters as special cases). Finally, Section 15.6 examines the problems faced when keeping track of more than one thing.

## 15.1 TIME AND UNCERTAINTY

---

We have developed our techniques for probabilistic reasoning in the context of *static* worlds, in which each random variable has a single fixed value. For example, when repairing a car, we assume that whatever is broken remains broken during the process of diagnosis; our job is to infer the state of the car from observed evidence, which also remains fixed.

Now consider a slightly different problem: treating a diabetic patient. As in the case of car repair, we have evidence such as recent insulin doses, food intake, blood sugar measurements, and other physical signs. The task is to assess the current state of the patient, including the actual blood sugar level and insulin level. Given this information, we can make a decision about the patient's food intake and insulin dose. Unlike the case of car repair, here the *dynamic* aspects of the problem are essential. Blood sugar levels and measurements thereof can change rapidly over time, depending on recent food intake and insulin doses, metabolic activity, the time of day, and so on. To assess the current state from the history of evidence and to predict the outcomes of treatment actions, we must model these changes.

The same considerations arise in many other contexts, such as tracking the location of a robot, tracking the economic activity of a nation, and making sense of a spoken or written sequence of words. How can dynamic situations like these be modeled?

### 15.1.1 States and observations

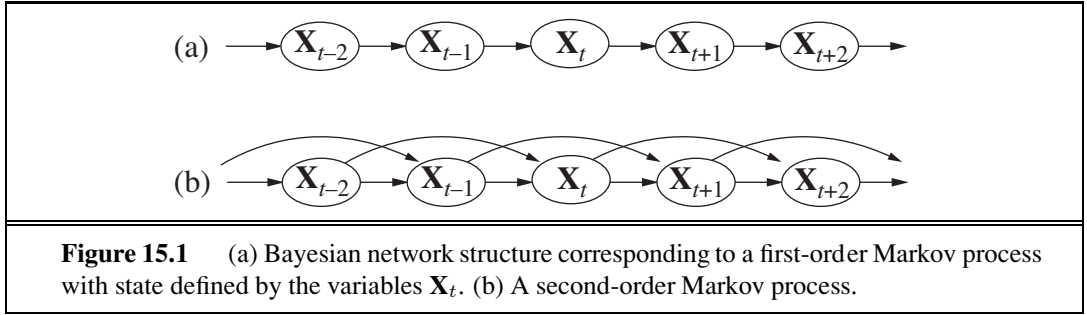
TIME SLICE

We view the world as a series of snapshots, or **time slices**, each of which contains a set of random variables, some observable and some not.<sup>1</sup> For simplicity, we will assume that the same subset of variables is observable in each time slice (although this is not strictly necessary in anything that follows). We will use  $\mathbf{X}_t$  to denote the set of state variables at time  $t$ , which are assumed to be unobservable, and  $\mathbf{E}_t$  to denote the set of observable evidence variables. The observation at time  $t$  is  $\mathbf{E}_t = \mathbf{e}_t$  for some set of values  $\mathbf{e}_t$ .

Consider the following example: You are the security guard stationed at a secret underground installation. You want to know whether it's raining today, but your only access to the outside world occurs each morning when you see the director coming in with, or without, an umbrella. For each day  $t$ , the set  $\mathbf{E}_t$  thus contains a single evidence variable  $Umbrella_t$  or  $U_t$  for short (whether the umbrella appears), and the set  $\mathbf{X}_t$  contains a single state variable  $Rain_t$  or  $R_t$  for short (whether it is raining). Other problems can involve larger sets of variables. In the diabetes example, we might have evidence variables, such as  $MeasuredBloodSugar_t$  and  $PulseRate_t$ , and state variables, such as  $BloodSugar_t$  and  $StomachContents_t$ . (Notice that  $BloodSugar_t$  and  $MeasuredBloodSugar_t$  are not the same variable; this is how we deal with noisy measurements of actual quantities.)

The interval between time slices also depends on the problem. For diabetes monitoring, a suitable interval might be an hour rather than a day. In this chapter we assume the interval between slices is fixed, so we can label times by integers. We will assume that the state sequence starts at  $t = 0$ ; for various uninteresting reasons, we will assume that evidence starts arriving at  $t = 1$  rather than  $t = 0$ . Hence, our umbrella world is represented by state variables  $R_0, R_1, R_2, \dots$  and evidence variables  $U_1, U_2, \dots$ . We will use the notation  $a:b$  to denote the sequence of integers from  $a$  to  $b$  (inclusive), and the notation  $\mathbf{X}_{a:b}$  to denote the set of variables from  $\mathbf{X}_a$  to  $\mathbf{X}_b$ . For example,  $U_{1:3}$  corresponds to the variables  $U_1, U_2, U_3$ .

<sup>1</sup> Uncertainty over *continuous* time can be modeled by **stochastic differential equations** (SDEs). The models studied in this chapter can be viewed as discrete-time approximations to SDEs.



### 15.1.2 Transition and sensor models

With the set of state and evidence variables for a given problem decided on, the next step is to specify how the world evolves (the transition model) and how the evidence variables get their values (the sensor model).

The transition model specifies the probability distribution over the latest state variables, given the previous values, that is,  $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{0:t-1})$ . Now we face a problem: the set  $\mathbf{X}_{0:t-1}$  is unbounded in size as  $t$  increases. We solve the problem by making a **Markov assumption**—that the current state depends on only a *finite fixed number* of previous states. Processes satisfying this assumption were first studied in depth by the Russian statistician Andrei Markov (1856–1922) and are called **Markov processes** or **Markov chains**. They come in various flavors; the simplest is the **first-order Markov process**, in which the current state depends only on the previous state and not on any earlier states. In other words, a state provides enough information to make the future conditionally independent of the past, and we have

$$\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{0:t-1}) = \mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-1}). \quad (15.1)$$

Hence, in a first-order Markov process, the transition model is the conditional distribution  $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-1})$ . The transition model for a second-order Markov process is the conditional distribution  $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-2}, \mathbf{X}_{t-1})$ . Figure 15.1 shows the Bayesian network structures corresponding to first-order and second-order Markov processes.

Even with the Markov assumption there is still a problem: there are infinitely many possible values of  $t$ . Do we need to specify a different distribution for each time step? We avoid this problem by assuming that changes in the world state are caused by a **stationary process**—that is, a process of change that is governed by laws that do not themselves change over time. (Don't confuse *stationary* with *static*: in a *static* process, the state itself does not change.) In the umbrella world, then, the conditional probability of rain,  $\mathbf{P}(R_t | R_{t-1})$ , is the same for all  $t$ , and we only have to specify one conditional probability table.

Now for the sensor model. The evidence variables  $\mathbf{E}_t$  *could* depend on previous variables as well as the current state variables, but any state that's worth its salt should suffice to generate the current sensor values. Thus, we make a **sensor Markov assumption** as follows:

$$\mathbf{P}(\mathbf{E}_t | \mathbf{X}_{0:t}, \mathbf{E}_{0:t-1}) = \mathbf{P}(\mathbf{E}_t | \mathbf{X}_t). \quad (15.2)$$

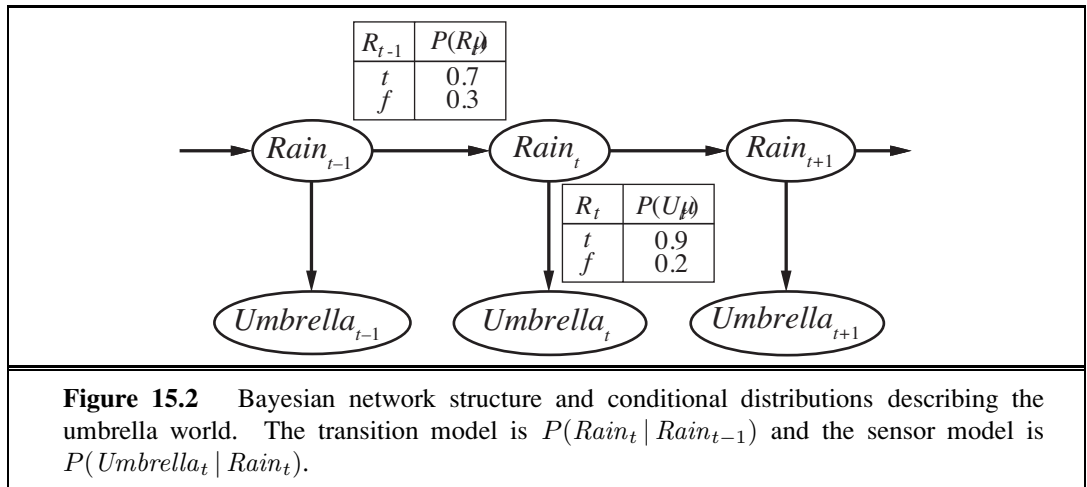
Thus,  $\mathbf{P}(\mathbf{E}_t | \mathbf{X}_t)$  is our sensor model (sometimes called the **observation model**). Figure 15.2 shows both the transition model and the sensor model for the umbrella example. Notice the

MARKOV  
ASSUMPTION

MARKOV PROCESS  
FIRST-ORDER  
MARKOV PROCESS

STATIONARY  
PROCESS

SENSOR MARKOV  
ASSUMPTION



direction of the dependence between state and sensors: the arrows go from the actual state of the world to sensor values because the state of the world *causes* the sensors to take on particular values: the rain *causes* the umbrella to appear. (The inference process, of course, goes in the other direction; the distinction between the direction of modeled dependencies and the direction of inference is one of the principal advantages of Bayesian networks.)

In addition to specifying the transition and sensor models, we need to say how everything gets started—the prior probability distribution at time 0,  $\mathbf{P}(\mathbf{X}_0)$ . With that, we have a specification of the complete joint distribution over all the variables, using Equation (14.2). For any  $t$ ,

$$\mathbf{P}(\mathbf{X}_{0:t}, \mathbf{E}_{1:t}) = \mathbf{P}(\mathbf{X}_0) \prod_{i=1}^t \mathbf{P}(\mathbf{X}_i | \mathbf{X}_{i-1}) \mathbf{P}(\mathbf{E}_i | \mathbf{X}_i). \quad (15.3)$$

The three terms on the right-hand side are the initial state model  $\mathbf{P}(\mathbf{X}_0)$ , the transition model  $\mathbf{P}(\mathbf{X}_i | \mathbf{X}_{i-1})$ , and the sensor model  $\mathbf{P}(\mathbf{E}_i | \mathbf{X}_i)$ .

The structure in Figure 15.2 is a first-order Markov process—the probability of rain is assumed to depend only on whether it rained the previous day. Whether such an assumption is reasonable depends on the domain itself. The first-order Markov assumption says that the state variables contain *all* the information needed to characterize the probability distribution for the next time slice. Sometimes the assumption is exactly true—for example, if a particle is executing a random walk along the  $x$ -axis, changing its position by  $\pm 1$  at each time step, then using the  $x$ -coordinate as the state gives a first-order Markov process. Sometimes the assumption is only approximate, as in the case of predicting rain only on the basis of whether it rained the previous day. There are two ways to improve the accuracy of the approximation:

1. Increasing the order of the Markov process model. For example, we could make a second-order model by adding  $Rain_{t-2}$  as a parent of  $Rain_t$ , which might give slightly more accurate predictions. For example, in Palo Alto, California, it very rarely rains more than two days in a row.
2. Increasing the set of state variables. For example, we could add  $Season_t$  to allow

us to incorporate historical records of rainy seasons, or we could add  $Temperature_t$ ,  $Humidity_t$  and  $Pressure_t$  (perhaps at a range of locations) to allow us to use a physical model of rainy conditions.

Exercise 15.1 asks you to show that the first solution—increasing the order—can always be reformulated as an increase in the set of state variables, keeping the order fixed. Notice that adding state variables might improve the system’s predictive power but also increases the prediction *requirements*: we now have to predict the new variables as well. Thus, we are looking for a “self-sufficient” set of variables, which really means that we have to understand the “physics” of the process being modeled. The requirement for accurate modeling of the process is obviously lessened if we can add new sensors (e.g., measurements of temperature and pressure) that provide information directly about the new state variables.

Consider, for example, the problem of tracking a robot wandering randomly on the X–Y plane. One might propose that the position and velocity are a sufficient set of state variables: one can simply use Newton’s laws to calculate the new position, and the velocity may change unpredictably. If the robot is battery-powered, however, then battery exhaustion would tend to have a systematic effect on the change in velocity. Because this in turn depends on how much power was used by all previous maneuvers, the Markov property is violated. We can restore the Markov property by including the charge level  $Battery_t$  as one of the state variables that make up  $\mathbf{X}_t$ . This helps in predicting the motion of the robot, but in turn requires a model for predicting  $Battery_t$  from  $Battery_{t-1}$  and the velocity. In some cases, that can be done reliably, but more often we find that error accumulates over time. In that case, accuracy can be improved by *adding a new sensor* for the battery level.

## 15.2 INFERENCE IN TEMPORAL MODELS

Having set up the structure of a generic temporal model, we can formulate the basic inference tasks that must be solved:

FILTERING  
BELIEF STATE  
STATE ESTIMATION

- **Filtering:** This is the task of computing the **belief state**—the posterior distribution over the most recent state—given all evidence to date. Filtering<sup>2</sup> is also called **state estimation**. In our example, we wish to compute  $\mathbf{P}(\mathbf{X}_t | \mathbf{e}_{1:t})$ . In the umbrella example, this would mean computing the probability of rain today, given all the observations of the umbrella carrier made so far. Filtering is what a rational agent does to keep track of the current state so that rational decisions can be made. It turns out that an almost identical calculation provides the likelihood of the evidence sequence,  $P(\mathbf{e}_{1:t})$ .

PREDICTION

- **Prediction:** This is the task of computing the posterior distribution over the *future* state, given all evidence to date. That is, we wish to compute  $\mathbf{P}(\mathbf{X}_{t+k} | \mathbf{e}_{1:t})$  for some  $k > 0$ . In the umbrella example, this might mean computing the probability of rain three days from now, given all the observations to date. Prediction is useful for evaluating possible courses of action based on their expected outcomes.

<sup>2</sup> The term “filtering” refers to the roots of this problem in early work on signal processing, where the problem is to filter out the noise in a signal by estimating its underlying properties.

## SMOOTHING

- **Smoothing:** This is the task of computing the posterior distribution over a *past* state, given all evidence up to the present. That is, we wish to compute  $\mathbf{P}(\mathbf{X}_k | \mathbf{e}_{1:t})$  for some  $k$  such that  $0 \leq k < t$ . In the umbrella example, it might mean computing the probability that it rained last Wednesday, given all the observations of the umbrella carrier made up to today. Smoothing provides a better estimate of the state than was available at the time, because it incorporates more evidence.<sup>3</sup>
- **Most likely explanation:** Given a sequence of observations, we might wish to find the sequence of states that is most likely to have generated those observations. That is, we wish to compute  $\text{argmax}_{\mathbf{x}_{1:t}} P(\mathbf{x}_{1:t} | \mathbf{e}_{1:t})$ . For example, if the umbrella appears on each of the first three days and is absent on the fourth, then the most likely explanation is that it rained on the first three days and did not rain on the fourth. Algorithms for this task are useful in many applications, including speech recognition—where the aim is to find the most likely sequence of words, given a series of sounds—and the reconstruction of bit strings transmitted over a noisy channel.

In addition to these inference tasks, we also have

- **Learning:** The transition and sensor models, if not yet known, can be learned from observations. Just as with static Bayesian networks, dynamic Bayes net learning can be done as a by-product of inference. Inference provides an estimate of what transitions actually occurred and of what states generated the sensor readings, and these estimates can be used to update the models. The updated models provide new estimates, and the process iterates to convergence. The overall process is an instance of the expectation-maximization or **EM algorithm**. (See Section 20.3.)

Note that learning requires smoothing, rather than filtering, because smoothing provides better estimates of the states of the process. Learning with filtering can fail to converge correctly; consider, for example, the problem of learning to solve murders: unless you are an eyewitness, smoothing is *always* required to infer what happened at the murder scene from the observable variables.

The remainder of this section describes generic algorithms for the four inference tasks, independent of the particular kind of model employed. Improvements specific to each model are described in subsequent sections.

### 15.2.1 Filtering and prediction

As we pointed out in Section 7.7.3, a useful filtering algorithm needs to maintain a current state estimate and update it, rather than going back over the entire history of percepts for each update. (Otherwise, the cost of each update increases as time goes by.) In other words, given the result of filtering up to time  $t$ , the agent needs to compute the result for  $t + 1$  from the new evidence  $\mathbf{e}_{t+1}$ ,

$$\mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) = f(\mathbf{e}_{t+1}, \mathbf{P}(\mathbf{X}_t | \mathbf{e}_{1:t})) ,$$

for some function  $f$ . This process is called **recursive estimation**. We can view the calculation

<sup>3</sup> In particular, when tracking a moving object with inaccurate position observations, smoothing gives a smoother estimated trajectory than filtering—hence the name.

as being composed of two parts: first, the current state distribution is projected forward from  $t$  to  $t+1$ ; then it is updated using the new evidence  $\mathbf{e}_{t+1}$ . This two-part process emerges quite simply when the formula is rearranged:

$$\begin{aligned} \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) &= \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}, \mathbf{e}_{t+1}) \quad (\text{dividing up the evidence}) \\ &= \alpha \mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}, \mathbf{e}_{1:t}) \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}) \quad (\text{using Bayes' rule}) \\ &= \alpha \mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}) \quad (\text{by the sensor Markov assumption}). \end{aligned} \quad (15.4)$$

Here and throughout this chapter,  $\alpha$  is a normalizing constant used to make probabilities sum up to 1. The second term,  $\mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t})$  represents a one-step prediction of the next state, and the first term updates this with the new evidence; notice that  $\mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1})$  is obtainable directly from the sensor model. Now we obtain the one-step prediction for the next state by conditioning on the current state  $\mathbf{X}_t$ :

$$\begin{aligned} \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) &= \alpha \mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{x}_t, \mathbf{e}_{1:t}) P(\mathbf{x}_t | \mathbf{e}_{1:t}) \\ &= \alpha \mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t}) \quad (\text{Markov assumption}). \end{aligned} \quad (15.5)$$

Within the summation, the first factor comes from the transition model and the second comes from the current state distribution. Hence, we have the desired recursive formulation. We can think of the filtered estimate  $\mathbf{P}(\mathbf{X}_t | \mathbf{e}_{1:t})$  as a “message”  $\mathbf{f}_{1:t}$  that is propagated forward along the sequence, modified by each transition and updated by each new observation. The process is given by

$$\mathbf{f}_{1:t+1} = \alpha \text{FORWARD}(\mathbf{f}_{1:t}, \mathbf{e}_{t+1}),$$

where FORWARD implements the update described in Equation (15.5) and the process begins with  $\mathbf{f}_{1:0} = \mathbf{P}(\mathbf{X}_0)$ . When all the state variables are discrete, the time for each update is constant (i.e., independent of  $t$ ), and the space required is also constant. (The constants depend, of course, on the size of the state space and the specific type of the temporal model in question.) *The time and space requirements for updating must be constant if an agent with limited memory is to keep track of the current state distribution over an unbounded sequence of observations.*

Let us illustrate the filtering process for two steps in the basic umbrella example (Figure 15.2.) That is, we will compute  $\mathbf{P}(R_2 | u_{1:2})$  as follows:

- On day 0, we have no observations, only the security guard’s prior beliefs; let’s assume that consists of  $\mathbf{P}(R_0) = \langle 0.5, 0.5 \rangle$ .
- On day 1, the umbrella appears, so  $U_1 = \text{true}$ . The prediction from  $t = 0$  to  $t = 1$  is

$$\begin{aligned} \mathbf{P}(R_1) &= \sum_{r_0} \mathbf{P}(R_1 | r_0) P(r_0) \\ &= \langle 0.7, 0.3 \rangle \times 0.5 + \langle 0.3, 0.7 \rangle \times 0.5 = \langle 0.5, 0.5 \rangle. \end{aligned}$$

Then the update step simply multiplies by the probability of the evidence for  $t = 1$  and normalizes, as shown in Equation (15.4):

$$\begin{aligned} \mathbf{P}(R_1 | u_1) &= \alpha \mathbf{P}(u_1 | R_1) \mathbf{P}(R_1) = \alpha \langle 0.9, 0.2 \rangle \langle 0.5, 0.5 \rangle \\ &= \alpha \langle 0.45, 0.1 \rangle \approx \langle 0.818, 0.182 \rangle. \end{aligned}$$





- On day 2, the umbrella appears, so  $U_2 = \text{true}$ . The prediction from  $t = 1$  to  $t = 2$  is

$$\begin{aligned} \mathbf{P}(R_2 | u_1) &= \sum_{r_1} \mathbf{P}(R_2 | r_1) P(r_1 | u_1) \\ &= \langle 0.7, 0.3 \rangle \times 0.818 + \langle 0.3, 0.7 \rangle \times 0.182 \approx \langle 0.627, 0.373 \rangle, \end{aligned}$$

and updating it with the evidence for  $t = 2$  gives

$$\begin{aligned} \mathbf{P}(R_2 | u_1, u_2) &= \alpha \mathbf{P}(u_2 | R_2) \mathbf{P}(R_2 | u_1) = \alpha \langle 0.9, 0.2 \rangle \langle 0.627, 0.373 \rangle \\ &= \alpha \langle 0.565, 0.075 \rangle \approx \langle 0.883, 0.117 \rangle. \end{aligned}$$

Intuitively, the probability of rain increases from day 1 to day 2 because rain persists. Exercise 15.2(a) asks you to investigate this tendency further.

The task of **prediction** can be seen simply as filtering without the addition of new evidence. In fact, the filtering process already incorporates a one-step prediction, and it is easy to derive the following recursive computation for predicting the state at  $t + k + 1$  from a prediction for  $t + k$ :

$$\mathbf{P}(\mathbf{X}_{t+k+1} | \mathbf{e}_{1:t}) = \sum_{\mathbf{x}_{t+k}} \mathbf{P}(\mathbf{X}_{t+k+1} | \mathbf{x}_{t+k}) P(\mathbf{x}_{t+k} | \mathbf{e}_{1:t}). \quad (15.6)$$

Naturally, this computation involves only the transition model and not the sensor model.

It is interesting to consider what happens as we try to predict further and further into the future. As Exercise 15.2(b) shows, the predicted distribution for rain converges to a fixed point  $\langle 0.5, 0.5 \rangle$ , after which it remains constant for all time. This is the **stationary distribution** of the Markov process defined by the transition model. (See also page 537.) A great deal is known about the properties of such distributions and about the **mixing time**—roughly, the time taken to reach the fixed point. In practical terms, this dooms to failure any attempt to predict the *actual* state for a number of steps that is more than a small fraction of the mixing time, unless the stationary distribution itself is strongly peaked in a small area of the state space. The more uncertainty there is in the transition model, the shorter will be the mixing time and the more the future is obscured.

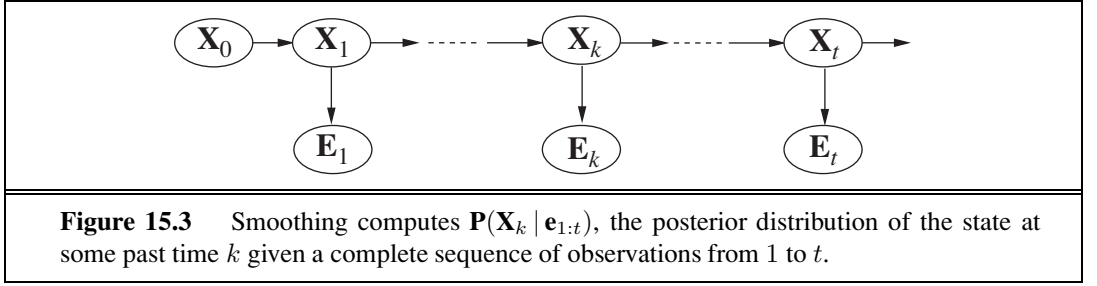
In addition to filtering and prediction, we can use a forward recursion to compute the **likelihood** of the evidence sequence,  $P(\mathbf{e}_{1:t})$ . This is a useful quantity if we want to compare different temporal models that might have produced the same evidence sequence (e.g., two different models for the persistence of rain). For this recursion, we use a likelihood message  $\ell_{1:t}(\mathbf{X}_t) = \mathbf{P}(\mathbf{X}_t, \mathbf{e}_{1:t})$ . It is a simple exercise to show that the message calculation is identical to that for filtering:

$$\ell_{1:t+1} = \text{FORWARD}(\ell_{1:t}, \mathbf{e}_{t+1}).$$

Having computed  $\ell_{1:t}$ , we obtain the actual likelihood by summing out  $\mathbf{X}_t$ :

$$L_{1:t} = P(\mathbf{e}_{1:t}) = \sum_{\mathbf{x}_t} \ell_{1:t}(\mathbf{x}_t). \quad (15.7)$$

Notice that the likelihood message represents the probabilities of longer and longer evidence sequences as time goes by and so becomes numerically smaller and smaller, leading to underflow problems with floating-point arithmetic. This is an important problem in practice, but we shall not go into solutions here.



### 15.2.2 Smoothing

As we said earlier, smoothing is the process of computing the distribution over past states given evidence up to the present; that is,  $\mathbf{P}(\mathbf{X}_k | \mathbf{e}_{1:t})$  for  $0 \leq k < t$ . (See Figure 15.3.) In anticipation of another recursive message-passing approach, we can split the computation into two parts—the evidence up to  $k$  and the evidence from  $k + 1$  to  $t$ ,

$$\begin{aligned}
 \mathbf{P}(\mathbf{X}_k | \mathbf{e}_{1:t}) &= \mathbf{P}(\mathbf{X}_k | \mathbf{e}_{1:k}, \mathbf{e}_{k+1:t}) \\
 &= \alpha \mathbf{P}(\mathbf{X}_k | \mathbf{e}_{1:k}) \mathbf{P}(\mathbf{e}_{k+1:t} | \mathbf{X}_k, \mathbf{e}_{1:k}) \quad (\text{using Bayes' rule}) \\
 &= \alpha \mathbf{P}(\mathbf{X}_k | \mathbf{e}_{1:k}) \mathbf{P}(\mathbf{e}_{k+1:t} | \mathbf{X}_k) \quad (\text{using conditional independence}) \\
 &= \alpha \mathbf{f}_{1:k} \times \mathbf{b}_{k+1:t} .
 \end{aligned} \tag{15.8}$$

where “ $\times$ ” represents pointwise multiplication of vectors. Here we have defined a “backward” message  $\mathbf{b}_{k+1:t} = \mathbf{P}(\mathbf{e}_{k+1:t} | \mathbf{X}_k)$ , analogous to the forward message  $\mathbf{f}_{1:k}$ . The forward message  $\mathbf{f}_{1:k}$  can be computed by filtering forward from 1 to  $k$ , as given by Equation (15.5). It turns out that the backward message  $\mathbf{b}_{k+1:t}$  can be computed by a recursive process that runs *backward* from  $t$ :

$$\begin{aligned}
 \mathbf{P}(\mathbf{e}_{k+1:t} | \mathbf{X}_k) &= \sum_{\mathbf{x}_{k+1}} \mathbf{P}(\mathbf{e}_{k+1:t} | \mathbf{X}_k, \mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) \quad (\text{conditioning on } \mathbf{X}_{k+1}) \\
 &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1:t} | \mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) \quad (\text{by conditional independence}) \\
 &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1}, \mathbf{e}_{k+2:t} | \mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) \\
 &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1} | \mathbf{x}_{k+1}) P(\mathbf{e}_{k+2:t} | \mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) ,
 \end{aligned} \tag{15.9}$$

where the last step follows by the conditional independence of  $\mathbf{e}_{k+1}$  and  $\mathbf{e}_{k+2:t}$ , given  $\mathbf{X}_{k+1}$ . Of the three factors in this summation, the first and third are obtained directly from the model, and the second is the “recursive call.” Using the message notation, we have

$$\mathbf{b}_{k+1:t} = \text{BACKWARD}(\mathbf{b}_{k+2:t}, \mathbf{e}_{k+1}) ,$$

where BACKWARD implements the update described in Equation (15.9). As with the forward recursion, the time and space needed for each update are constant and thus independent of  $t$ .

We can now see that the two terms in Equation (15.8) can both be computed by recursions through time, one running forward from 1 to  $k$  and using the filtering equation (15.5)

and the other running backward from  $t$  to  $k + 1$  and using Equation (15.9). Note that the backward phase is initialized with  $\mathbf{b}_{t+1:t} = \mathbf{P}(\mathbf{e}_{t+1:t} | \mathbf{X}_t) = \mathbf{P}(\mathbf{e}_{t+1:t} | \mathbf{X}_t) \mathbf{1}$ , where  $\mathbf{1}$  is a vector of 1s. (Because  $\mathbf{e}_{t+1:t}$  is an empty sequence, the probability of observing it is 1.)

Let us now apply this algorithm to the umbrella example, computing the smoothed estimate for the probability of rain at time  $k = 1$ , given the umbrella observations on days 1 and 2. From Equation (15.8), this is given by

$$\mathbf{P}(R_1 | u_1, u_2) = \alpha \mathbf{P}(R_1 | u_1) \mathbf{P}(u_2 | R_1). \quad (15.10)$$

The first term we already know to be  $\langle .818, .182 \rangle$ , from the forward filtering process described earlier. The second term can be computed by applying the backward recursion in Equation (15.9):

$$\begin{aligned} \mathbf{P}(u_2 | R_1) &= \sum_{r_2} P(u_2 | r_2) P(r_2 | R_1) \\ &= (0.9 \times 1 \times \langle 0.7, 0.3 \rangle) + (0.2 \times 1 \times \langle 0.3, 0.7 \rangle) = \langle 0.69, 0.41 \rangle. \end{aligned}$$

Plugging this into Equation (15.10), we find that the smoothed estimate for rain on day 1 is

$$\mathbf{P}(R_1 | u_1, u_2) = \alpha \langle 0.818, 0.182 \rangle \times \langle 0.69, 0.41 \rangle \approx \langle 0.883, 0.117 \rangle.$$

Thus, the smoothed estimate for rain on day 1 is *higher* than the filtered estimate (0.818) in this case. This is because the umbrella on day 2 makes it more likely to have rained on day 2; in turn, because rain tends to persist, that makes it more likely to have rained on day 1.

Both the forward and backward recursions take a constant amount of time per step; hence, the time complexity of smoothing with respect to evidence  $\mathbf{e}_{1:t}$  is  $O(t)$ . This is the complexity for smoothing at a particular time step  $k$ . If we want to smooth the whole sequence, one obvious method is simply to run the whole smoothing process once for each time step to be smoothed. This results in a time complexity of  $O(t^2)$ . A better approach uses a simple application of dynamic programming to reduce the complexity to  $O(t)$ . A clue appears in the preceding analysis of the umbrella example, where we were able to reuse the results of the forward-filtering phase. The key to the linear-time algorithm is to *record the results* of forward filtering over the whole sequence. Then we run the backward recursion from  $t$  down to 1, computing the smoothed estimate at each step  $k$  from the computed backward message  $\mathbf{b}_{k+1:t}$  and the stored forward message  $\mathbf{f}_{1:k}$ . The algorithm, aptly called the **forward-backward algorithm**, is shown in Figure 15.4.

The alert reader will have spotted that the Bayesian network structure shown in Figure 15.3 is a *polytree* as defined on page 528. This means that a straightforward application of the clustering algorithm also yields a linear-time algorithm that computes smoothed estimates for the entire sequence. It is now understood that the forward-backward algorithm is in fact a special case of the polytree propagation algorithm used with clustering methods (although the two were developed independently).

The forward-backward algorithm forms the computational backbone for many applications that deal with sequences of noisy observations. As described so far, it has two practical drawbacks. The first is that its space complexity can be too high when the state space is large and the sequences are long. It uses  $O(|\mathbf{f}|t)$  space where  $|\mathbf{f}|$  is the size of the representation of the forward message. The space requirement can be reduced to  $O(|\mathbf{f}| \log t)$  with a concomi-

tant increase in the time complexity by a factor of  $\log t$ , as shown in Exercise 15.3. In some cases (see Section 15.3), a constant-space algorithm can be used.

The second drawback of the basic algorithm is that it needs to be modified to work in an *online* setting where smoothed estimates must be computed for earlier time slices as new observations are continuously added to the end of the sequence. The most common requirement is for **fixed-lag smoothing**, which requires computing the smoothed estimate  $\mathbf{P}(\mathbf{X}_{t-d} | \mathbf{e}_{1:t})$  for fixed  $d$ . That is, smoothing is done for the time slice  $d$  steps behind the current time  $t$ ; as  $t$  increases, the smoothing has to keep up. Obviously, we can run the forward–backward algorithm over the  $d$ -step “window” as each new observation is added, but this seems inefficient. In Section 15.3, we will see that fixed-lag smoothing can, in some cases, be done in constant time per update, independent of the lag  $d$ .

FIXED-LAG  
SMOOTHING

### 15.2.3 Finding the most likely sequence

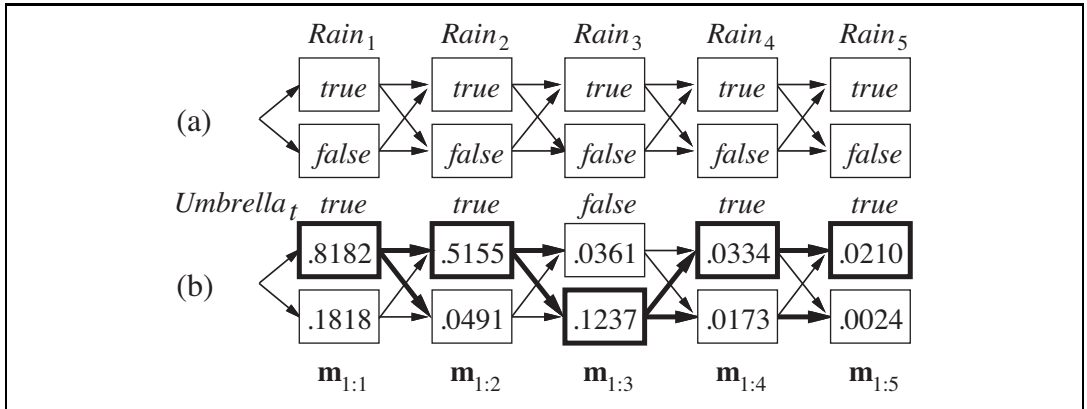
Suppose that  $[true, true, false, true, true]$  is the umbrella sequence for the security guard’s first five days on the job. What is the weather sequence most likely to explain this? Does the absence of the umbrella on day 3 mean that it wasn’t raining, or did the director forget to bring it? If it didn’t rain on day 3, perhaps (because weather tends to persist) it didn’t rain on day 4 either, but the director brought the umbrella just in case. In all, there are  $2^5$  possible weather sequences we could pick. Is there a way to find the most likely one, short of enumerating all of them?

We could try this linear-time procedure: use smoothing to find the posterior distribution for the weather at each time step; then construct the sequence, using at each step the weather that is most likely according to the posterior. Such an approach should set off alarm bells in the reader’s head, because the posterior distributions computed by smoothing are distri-

```
function FORWARD-BACKWARD(ev, prior) returns a vector of probability distributions
  inputs: ev, a vector of evidence values for steps 1, ...,  $t$ 
           prior, the prior distribution on the initial state,  $\mathbf{P}(\mathbf{X}_0)$ 
  local variables: fv, a vector of forward messages for steps 0, ...,  $t$ 
                   b, a representation of the backward message, initially all 1s
                   sv, a vector of smoothed estimates for steps 1, ...,  $t$ 

  fv[0]  $\leftarrow$  prior
  for  $i = 1$  to  $t$  do
    fv[ $i$ ]  $\leftarrow$  FORWARD(fv[ $i - 1$ ], ev[ $i$ ])
  for  $i = t$  downto 1 do
    sv[ $i$ ]  $\leftarrow$  NORMALIZE(fv[ $i$ ]  $\times$  b)
    b  $\leftarrow$  BACKWARD(b, ev[ $i$ ])
  return sv
```

**Figure 15.4** The forward–backward algorithm for smoothing: computing posterior probabilities of a sequence of states given a sequence of observations. The FORWARD and BACKWARD operators are defined by Equations (15.5) and (15.9), respectively.



**Figure 15.5** (a) Possible state sequences for  $Rain_t$  can be viewed as paths through a graph of the possible states at each time step. (States are shown as rectangles to avoid confusion with nodes in a Bayes net.) (b) Operation of the Viterbi algorithm for the umbrella observation sequence  $[true, true, false, true, true]$ . For each  $t$ , we have shown the values of the message  $m_{1:t}$ , which gives the probability of the best sequence reaching each state at time  $t$ . Also, for each state, the bold arrow leading into it indicates its best predecessor as measured by the product of the preceding sequence probability and the transition probability. Following the bold arrows back from the most likely state in  $m_{1:5}$  gives the most likely sequence.

butions over *single* time steps, whereas to find the most likely *sequence* we must consider *joint* probabilities over all the time steps. The results can in fact be quite different. (See Exercise 15.4.)

There *is* a linear-time algorithm for finding the most likely sequence, but it requires a little more thought. It relies on the same Markov property that yielded efficient algorithms for filtering and smoothing. The easiest way to think about the problem is to view each sequence as a *path* through a graph whose nodes are the possible *states* at each time step. Such a graph is shown for the umbrella world in Figure 15.5(a). Now consider the task of finding the most likely path through this graph, where the likelihood of any path is the product of the transition probabilities along the path and the probabilities of the given observations at each state. Let's focus in particular on paths that reach the state  $Rain_5 = true$ . Because of the Markov property, it follows that the most likely path to the state  $Rain_5 = true$  consists of the most likely path to *some* state at time 4 followed by a transition to  $Rain_5 = true$ ; and the state at time 4 that will become part of the path to  $Rain_5 = true$  is whichever maximizes the likelihood of that path. In other words, *there is a recursive relationship between most likely paths to each state  $x_{t+1}$  and most likely paths to each state  $x_t$* . We can write this relationship as an equation connecting the probabilities of the paths:

$$\begin{aligned} \max_{x_1 \dots x_t} P(x_1, \dots, x_t, x_{t+1} \mid e_{1:t+1}) \\ = \alpha P(e_{t+1} \mid x_{t+1}) \max_{x_t} \left( P(x_{t+1} \mid x_t) \max_{x_1 \dots x_{t-1}} P(x_1, \dots, x_{t-1}, x_t \mid e_{1:t}) \right). \end{aligned} \quad (15.11)$$

Equation (15.11) is *identical* to the filtering equation (15.5) except that

1. The forward message  $\mathbf{f}_{1:t} = \mathbf{P}(\mathbf{X}_t | \mathbf{e}_{1:t})$  is replaced by the message

$$\mathbf{m}_{1:t} = \max_{\mathbf{x}_1 \dots \mathbf{x}_{t-1}} \mathbf{P}(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}, \mathbf{X}_t | \mathbf{e}_{1:t}),$$

that is, the probabilities of the most likely path to each state  $\mathbf{x}_t$ ; and

2. the summation over  $\mathbf{x}_t$  in Equation (15.5) is replaced by the maximization over  $\mathbf{x}_t$  in Equation (15.11).

Thus, the algorithm for computing the most likely sequence is similar to filtering: it runs forward along the sequence, computing the  $\mathbf{m}$  message at each time step, using Equation (15.11). The progress of this computation is shown in Figure 15.5(b). At the end, it will have the probability for the most likely sequence reaching *each* of the final states. One can thus easily select the most likely sequence overall (the states outlined in bold). In order to identify the actual sequence, as opposed to just computing its probability, the algorithm will also need to record, for each state, the best state that leads to it; these are indicated by the bold arrows in Figure 15.5(b). The optimal sequence is identified by following these bold arrows backwards from the best final state.

VITERBI ALGORITHM

The algorithm we have just described is called the **Viterbi algorithm**, after its inventor. Like the filtering algorithm, its time complexity is linear in  $t$ , the length of the sequence. Unlike filtering, which uses constant space, its space requirement is also linear in  $t$ . This is because the Viterbi algorithm needs to keep the pointers that identify the best sequence leading to each state.

### 15.3 HIDDEN MARKOV MODELS

The preceding section developed algorithms for temporal probabilistic reasoning using a general framework that was independent of the specific form of the transition and sensor models. In this and the next two sections, we discuss more concrete models and applications that illustrate the power of the basic algorithms and in some cases allow further improvements.

HIDDEN MARKOV MODEL

We begin with the **hidden Markov model**, or **HMM**. An HMM is a temporal probabilistic model in which the state of the process is described by a *single discrete* random variable. The possible values of the variable are the possible states of the world. The umbrella example described in the preceding section is therefore an HMM, since it has just one state variable:  $Rain_t$ . What happens if you have a model with two or more state variables? You can still fit it into the HMM framework by combining the variables into a single “megavariable” whose values are all possible tuples of values of the individual state variables. We will see that the restricted structure of HMMs allows for a simple and elegant matrix implementation of all the basic algorithms.<sup>4</sup>

<sup>4</sup> The reader unfamiliar with basic operations on vectors and matrices might wish to consult Appendix A before proceeding with this section.

### 15.3.1 Simplified matrix algorithms

With a single, discrete state variable  $X_t$ , we can give concrete form to the representations of the transition model, the sensor model, and the forward and backward messages. Let the state variable  $X_t$  have values denoted by integers  $1, \dots, S$ , where  $S$  is the number of possible states. The transition model  $\mathbf{P}(X_t | X_{t-1})$  becomes an  $S \times S$  matrix  $\mathbf{T}$ , where

$$\mathbf{T}_{ij} = P(X_t = j | X_{t-1} = i) .$$

That is,  $\mathbf{T}_{ij}$  is the probability of a transition from state  $i$  to state  $j$ . For example, the transition matrix for the umbrella world is

$$\mathbf{T} = \mathbf{P}(X_t | X_{t-1}) = \begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix} .$$

We also put the sensor model in matrix form. In this case, because the value of the evidence variable  $E_t$  is known at time  $t$  (call it  $e_t$ ), we need only specify, for each state, how likely it is that the state causes  $e_t$  to appear: we need  $P(e_t | X_t = i)$  for each state  $i$ . For mathematical convenience we place these values into an  $S \times S$  diagonal matrix,  $\mathbf{O}_t$  whose  $i$ th diagonal entry is  $P(e_t | X_t = i)$  and whose other entries are 0. For example, on day 1 in the umbrella world of Figure 15.5,  $U_1 = \text{true}$ , and on day 3,  $U_3 = \text{false}$ , so, from Figure 15.2, we have

$$\mathbf{O}_1 = \begin{pmatrix} 0.9 & 0 \\ 0 & 0.2 \end{pmatrix}; \quad \mathbf{O}_3 = \begin{pmatrix} 0.1 & 0 \\ 0 & 0.8 \end{pmatrix} .$$

Now, if we use column vectors to represent the forward and backward messages, all the computations become simple matrix–vector operations. The forward equation (15.5) becomes

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^\top \mathbf{f}_{1:t} \quad (15.12)$$

and the backward equation (15.9) becomes

$$\mathbf{b}_{k+1:t} = \mathbf{T} \mathbf{O}_{k+1} \mathbf{b}_{k+2:t} . \quad (15.13)$$

From these equations, we can see that the time complexity of the forward–backward algorithm (Figure 15.4) applied to a sequence of length  $t$  is  $O(S^2 t)$ , because each step requires multiplying an  $S$ -element vector by an  $S \times S$  matrix. The space requirement is  $O(St)$ , because the forward pass stores  $t$  vectors of size  $S$ .

Besides providing an elegant description of the filtering and smoothing algorithms for HMMs, the matrix formulation reveals opportunities for improved algorithms. The first is a simple variation on the forward–backward algorithm that allows smoothing to be carried out in *constant* space, independently of the length of the sequence. The idea is that smoothing for any particular time slice  $k$  requires the simultaneous presence of both the forward and backward messages,  $\mathbf{f}_{1:k}$  and  $\mathbf{b}_{k+1:t}$ , according to Equation (15.8). The forward–backward algorithm achieves this by storing the  $\mathbf{f}$ s computed on the forward pass so that they are available during the backward pass. Another way to achieve this is with a single pass that propagates both  $\mathbf{f}$  and  $\mathbf{b}$  in the same direction. For example, the “forward” message  $\mathbf{f}$  can be propagated backward if we manipulate Equation (15.12) to work in the other direction:

$$\mathbf{f}_{1:t} = \alpha' (\mathbf{T}^\top)^{-1} \mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} .$$

The modified smoothing algorithm works by first running the standard forward pass to compute  $\mathbf{f}_{t:t}$  (forgetting all the intermediate results) and then running the backward pass for both

```

function FIXED-LAG-SMOOTHING( $e_t, hmm, d$ ) returns a distribution over  $\mathbf{X}_{t-d}$ 
  inputs:  $e_t$ , the current evidence for time step  $t$ 
            $hmm$ , a hidden Markov model with  $S \times S$  transition matrix  $\mathbf{T}$ 
            $d$ , the length of the lag for smoothing
  persistent:  $t$ , the current time, initially 1
                 $\mathbf{f}$ , the forward message  $\mathbf{P}(X_t|e_{1:t})$ , initially  $hmm.PRIOR$ 
                 $\mathbf{B}$ , the  $d$ -step backward transformation matrix, initially the identity matrix
                 $e_{t-d:t}$ , double-ended list of evidence from  $t-d$  to  $t$ , initially empty
  local variables:  $\mathbf{O}_{t-d}, \mathbf{O}_t$ , diagonal matrices containing the sensor model information

  add  $e_t$  to the end of  $e_{t-d:t}$ 
   $\mathbf{O}_t \leftarrow$  diagonal matrix containing  $\mathbf{P}(e_t|X_t)$ 
  if  $t > d$  then
     $\mathbf{f} \leftarrow \text{FORWARD}(\mathbf{f}, e_t)$ 
    remove  $e_{t-d-1}$  from the beginning of  $e_{t-d:t}$ 
     $\mathbf{O}_{t-d} \leftarrow$  diagonal matrix containing  $\mathbf{P}(e_{t-d}|X_{t-d})$ 
     $\mathbf{B} \leftarrow \mathbf{O}_{t-d}^{-1} \mathbf{T}^{-1} \mathbf{B} \mathbf{O}_t$ 
  else  $\mathbf{B} \leftarrow \mathbf{B} \mathbf{O}_t$ 
   $t \leftarrow t + 1$ 
  if  $t > d$  then return  $\text{NORMALIZE}(\mathbf{f} \times \mathbf{B} \mathbf{1})$  else return null

```

**Figure 15.6** An algorithm for smoothing with a fixed time lag of  $d$  steps, implemented as an online algorithm that outputs the new smoothed estimate given the observation for a new time step. Notice that the final output  $\text{NORMALIZE}(\mathbf{f} \times \mathbf{B} \mathbf{1})$  is just  $\alpha \mathbf{f} \times \mathbf{b}$ , by Equation (15.14).

$\mathbf{b}$  and  $\mathbf{f}$  together, using them to compute the smoothed estimate at each step. Since only one copy of each message is needed, the storage requirements are constant (i.e., independent of  $t$ , the length of the sequence). There are two significant restrictions on this algorithm: it requires that the transition matrix be invertible and that the sensor model have no zeroes—that is, that every observation be possible in every state.

A second area in which the matrix formulation reveals an improvement is in *online* smoothing with a fixed lag. The fact that smoothing can be done in constant space suggests that there should exist an efficient recursive algorithm for online smoothing—that is, an algorithm whose time complexity is independent of the length of the lag. Let us suppose that the lag is  $d$ ; that is, we are smoothing at time slice  $t-d$ , where the current time is  $t$ . By Equation (15.8), we need to compute

$$\alpha \mathbf{f}_{1:t-d} \times \mathbf{b}_{t-d+1:t}$$

for slice  $t-d$ . Then, when a new observation arrives, we need to compute

$$\alpha \mathbf{f}_{1:t-d+1} \times \mathbf{b}_{t-d+2:t+1}$$

for slice  $t-d+1$ . How can this be done incrementally? First, we can compute  $\mathbf{f}_{1:t-d+1}$  from  $\mathbf{f}_{1:t-d}$ , using the standard filtering process, Equation (15.5).



Computing the backward message incrementally is trickier, because there is no simple relationship between the old backward message  $\mathbf{b}_{t-d+1:t}$  and the new backward message  $\mathbf{b}_{t-d+2:t+1}$ . Instead, we will examine the relationship between the old backward message  $\mathbf{b}_{t-d+1:t}$  and the backward message at the front of the sequence,  $\mathbf{b}_{t+1:t}$ . To do this, we apply Equation (15.13)  $d$  times to get

$$\mathbf{b}_{t-d+1:t} = \left( \prod_{i=t-d+1}^t \mathbf{TO}_i \right) \mathbf{b}_{t+1:t} = \mathbf{B}_{t-d+1:t} \mathbf{1}, \quad (15.14)$$

where the matrix  $\mathbf{B}_{t-d+1:t}$  is the product of the sequence of  $\mathbf{T}$  and  $\mathbf{O}$  matrices.  $\mathbf{B}$  can be thought of as a “transformation operator” that transforms a later backward message into an earlier one. A similar equation holds for the new backward messages *after* the next observation arrives:

$$\mathbf{b}_{t-d+2:t+1} = \left( \prod_{i=t-d+2}^{t+1} \mathbf{TO}_i \right) \mathbf{b}_{t+2:t+1} = \mathbf{B}_{t-d+2:t+1} \mathbf{1}. \quad (15.15)$$

Examining the product expressions in Equations (15.14) and (15.15), we see that they have a simple relationship: to get the second product, “divide” the first product by the first element  $\mathbf{TO}_{t-d+1}$ , and multiply by the new last element  $\mathbf{TO}_{t+1}$ . In matrix language, then, there is a simple relationship between the old and new  $\mathbf{B}$  matrices:

$$\mathbf{B}_{t-d+2:t+1} = \mathbf{O}_{t-d+1}^{-1} \mathbf{T}^{-1} \mathbf{B}_{t-d+1:t} \mathbf{TO}_{t+1}. \quad (15.16)$$

This equation provides an incremental update for the  $\mathbf{B}$  matrix, which in turn (through Equation (15.15)) allows us to compute the new backward message  $\mathbf{b}_{t-d+2:t+1}$ . The complete algorithm, which requires storing and updating  $\mathbf{f}$  and  $\mathbf{B}$ , is shown in Figure 15.6.

### 15.3.2 Hidden Markov model example: Localization

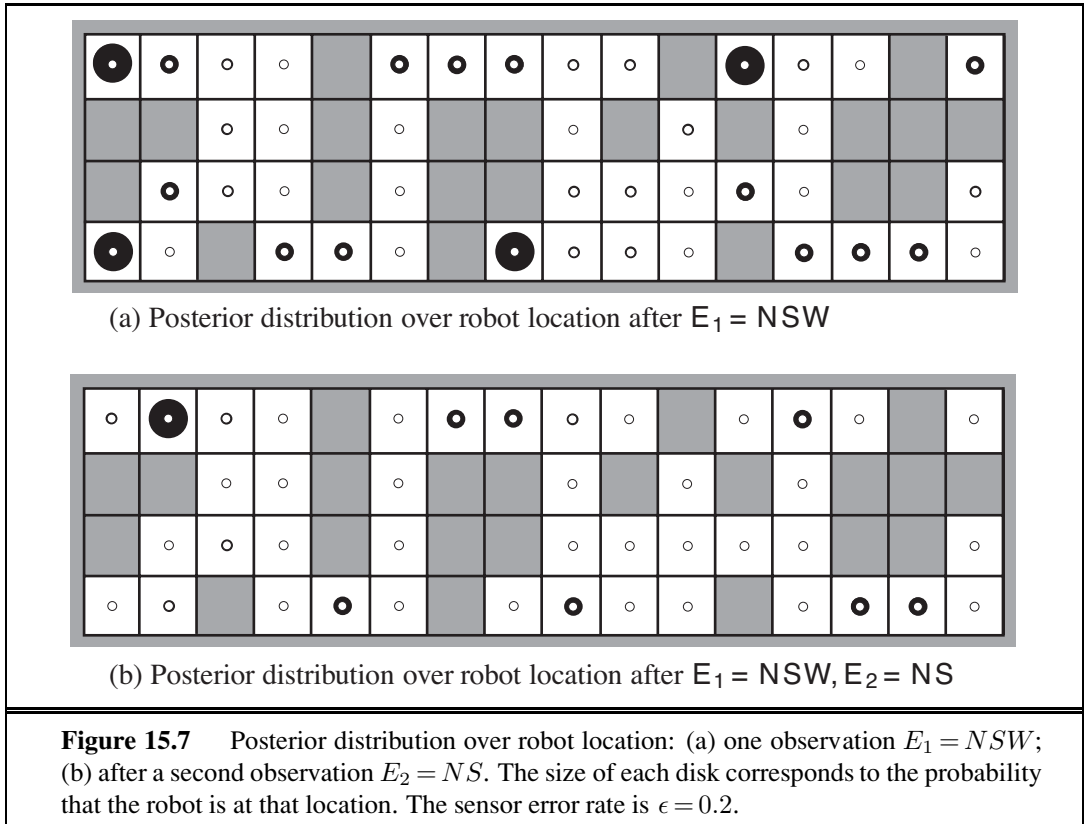
On page 145, we introduced a simple form of the **localization** problem for the vacuum world. In that version, the robot had a single nondeterministic *Move* action and its sensors reported perfectly whether or not obstacles lay immediately to the north, south, east, and west; the robot’s belief state was the set of possible locations it could be in.

Here we make the problem slightly more realistic by including a simple probability model for the robot’s motion and by allowing for noise in the sensors. The state variable  $X_t$  represents the location of the robot on the discrete grid; the domain of this variable is the set of empty squares  $\{s_1, \dots, s_n\}$ . Let  $\text{NEIGHBORS}(s)$  be the set of empty squares that are adjacent to  $s$  and let  $N(s)$  be the size of that set. Then the transition model for *Move* action says that the robot is equally likely to end up at any neighboring square:

$$P(X_{t+1} = j \mid X_t = i) = \mathbf{T}_{ij} = (1/N(i) \text{ if } j \in \text{NEIGHBORS}(i) \text{ else } 0).$$

We don’t know where the robot starts, so we will assume a uniform distribution over all the squares; that is,  $P(X_0 = i) = 1/n$ . For the particular environment we consider (Figure 15.7),  $n = 42$  and the transition matrix  $\mathbf{T}$  has  $42 \times 42 = 1764$  entries.

The sensor variable  $E_t$  has 16 possible values, each a four-bit sequence giving the presence or absence of an obstacle in a particular compass direction. We will use the notation



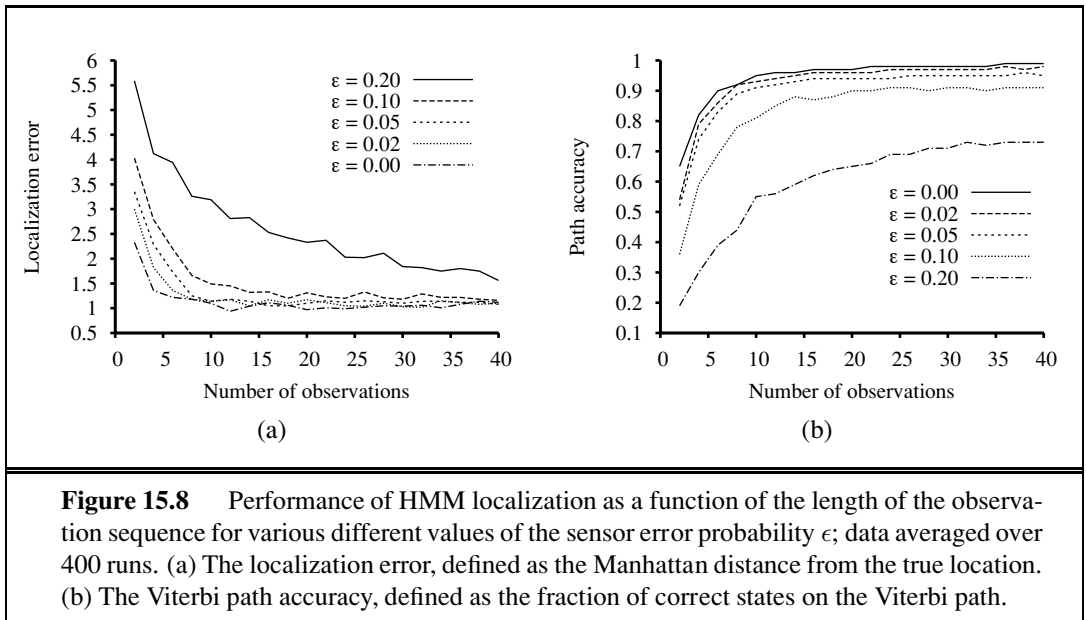
$\text{NS}$ , for example, to mean that the north and south sensors report an obstacle and the east and west do not. Suppose that each sensor's error rate is  $\epsilon$  and that errors occur independently for the four sensor directions. In that case, the probability of getting all four bits right is  $(1 - \epsilon)^4$  and the probability of getting them all wrong is  $\epsilon^4$ . Furthermore, if  $d_{it}$  is the discrepancy—the number of bits that are different—between the true values for square  $i$  and the actual reading  $e_t$ , then the probability that a robot in square  $i$  would receive a sensor reading  $e_t$  is

$$P(E_t = e_t \mid X_t = i) = \mathbf{O}_{t_{ii}} = (1 - \epsilon)^{4-d_{it}} \epsilon^{d_{it}}.$$

For example, the probability that a square with obstacles to the north and south would produce a sensor reading  $\text{NSE}$  is  $(1 - \epsilon)^3 \epsilon^1$ .

Given the matrices  $\mathbf{T}$  and  $\mathbf{O}_t$ , the robot can use Equation (15.12) to compute the posterior distribution over locations—that is, to work out where it is. Figure 15.7 shows the distributions  $\mathbf{P}(X_1 \mid E_1 = \text{NSW})$  and  $\mathbf{P}(X_2 \mid E_1 = \text{NSW}, E_2 = \text{NS})$ . This is the same maze we saw before in Figure 4.18 (page 146), but there we used logical filtering to find the locations that were *possible*, assuming perfect sensing. Those same locations are still the most *likely* with noisy sensing, but now *every* location has some nonzero probability.

In addition to filtering to estimate its current location, the robot can use smoothing (Equation (15.13)) to work out where it was at any given past time—for example, where it began at time 0—and it can use the Viterbi algorithm to work out the most likely path it has



taken to get where it is now. Figure 15.8 shows the localization error and Viterbi path accuracy for various values of the per-bit sensor error rate  $\epsilon$ . Even when  $\epsilon$  is 20%—which means that the overall sensor reading is wrong 59% of the time—the robot is usually able to work out its location within two squares after 25 observations. This is because of the algorithm’s ability to integrate evidence over time and to take into account the probabilistic constraints imposed on the location sequence by the transition model. When  $\epsilon$  is 10%, the performance after a half-dozen observations is hard to distinguish from the performance with perfect sensing. Exercise 15.7 asks you to explore how robust the HMM localization algorithm is to errors in the prior distribution  $\mathbf{P}(X_0)$  and in the transition model itself. Broadly speaking, high levels of localization and path accuracy are maintained even in the face of substantial errors in the models used.

The state variable for the example we have considered in this section is a physical location in the world. Other problems can, of course, include other aspects of the world. Exercise 15.8 asks you to consider a version of the vacuum robot that has the policy of going straight for as long as it can; only when it encounters an obstacle does it change to a new (randomly selected) heading. To model this robot, each state in the model consists of a (*location, heading*) pair. For the environment in Figure 15.7, which has 42 empty squares, this leads to  $168$  states and a transition matrix with  $168^2 = 28,224$  entries—still a manageable number. If we add the possibility of dirt in the squares, the number of states is multiplied by  $2^{42}$  and the transition matrix ends up with more than  $10^{29}$  entries—no longer a manageable number; Section 15.5 shows how to use dynamic Bayesian networks to model domains with many state variables. If we allow the robot to move continuously rather than in a discrete grid, the number of states becomes infinite; the next section shows how to handle this case.

## 15.4 KALMAN FILTERS

KALMAN FILTERING

Imagine watching a small bird flying through dense jungle foliage at dusk: you glimpse brief, intermittent flashes of motion; you try hard to guess where the bird is and where it will appear next so that you don't lose it. Or imagine that you are a World War II radar operator peering at a faint, wandering blip that appears once every 10 seconds on the screen. Or, going back further still, imagine you are Kepler trying to reconstruct the motions of the planets from a collection of highly inaccurate angular observations taken at irregular and imprecisely measured intervals. In all these cases, you are doing filtering: estimating state variables (here, position and velocity) from noisy observations over time. If the variables were discrete, we could model the system with a hidden Markov model. This section examines methods for handling continuous variables, using an algorithm called **Kalman filtering**, after one of its inventors, Rudolf E. Kalman.

The bird's flight might be specified by six continuous variables at each time point; three for position  $(X_t, Y_t, Z_t)$  and three for velocity  $(\dot{X}_t, \dot{Y}_t, \dot{Z}_t)$ . We will need suitable conditional densities to represent the transition and sensor models; as in Chapter 14, we will use **linear Gaussian** distributions. This means that the next state  $\mathbf{X}_{t+1}$  must be a linear function of the current state  $\mathbf{X}_t$ , plus some Gaussian noise, a condition that turns out to be quite reasonable in practice. Consider, for example, the  $X$ -coordinate of the bird, ignoring the other coordinates for now. Let the time interval between observations be  $\Delta$ , and assume constant velocity during the interval; then the position update is given by  $X_{t+\Delta} = X_t + \dot{X}_t \Delta$ . Adding Gaussian noise (to account for wind variation, etc.), we obtain a linear Gaussian transition model:

$$P(X_{t+\Delta} = x_{t+\Delta} \mid X_t = x_t, \dot{X}_t = \dot{x}_t) = N(x_t + \dot{x}_t \Delta, \sigma^2)(x_{t+\Delta}) .$$

The Bayesian network structure for a system with position vector  $\mathbf{X}_t$  and velocity  $\dot{\mathbf{X}}_t$  is shown in Figure 15.9. Note that this is a very specific form of linear Gaussian model; the general form will be described later in this section and covers a vast array of applications beyond the simple motion examples of the first paragraph. The reader might wish to consult Appendix A for some of the mathematical properties of Gaussian distributions; for our immediate purposes, the most important is that a **multivariate Gaussian** distribution for  $d$  variables is specified by a  $d$ -element mean  $\boldsymbol{\mu}$  and a  $d \times d$  covariance matrix  $\boldsymbol{\Sigma}$ .

MULTIVARIATE  
GAUSSIAN

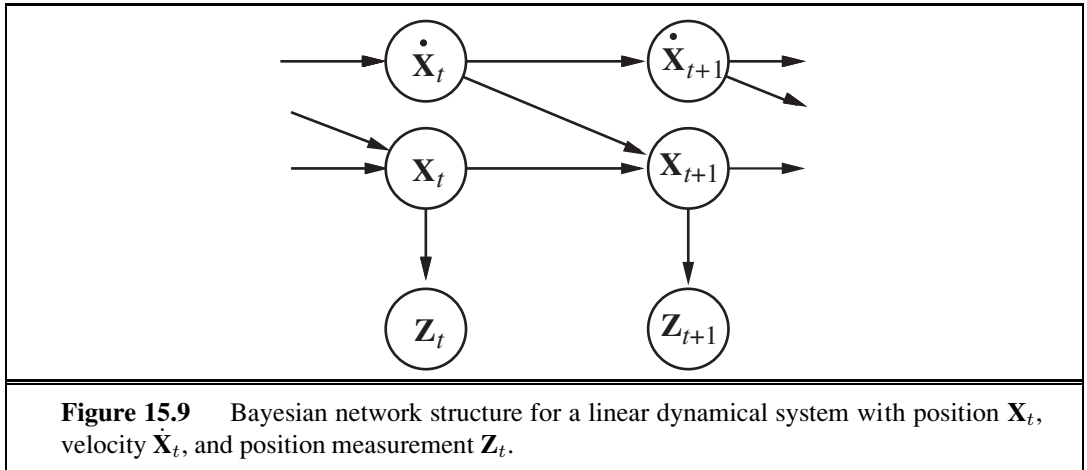
### 15.4.1 Updating Gaussian distributions

In Chapter 14 on page 521, we alluded to a key property of the linear Gaussian family of distributions: it remains closed under the standard Bayesian network operations. Here, we make this claim precise in the context of filtering in a temporal probability model. The required properties correspond to the two-step filtering calculation in Equation (15.5):

1. If the current distribution  $\mathbf{P}(\mathbf{X}_t \mid \mathbf{e}_{1:t})$  is Gaussian and the transition model  $\mathbf{P}(\mathbf{X}_{t+1} \mid \mathbf{x}_t)$  is linear Gaussian, then the one-step predicted distribution given by

$$\mathbf{P}(\mathbf{X}_{t+1} \mid \mathbf{e}_{1:t}) = \int_{\mathbf{x}_t} \mathbf{P}(\mathbf{X}_{t+1} \mid \mathbf{x}_t) P(\mathbf{x}_t \mid \mathbf{e}_{1:t}) d\mathbf{x}_t \quad (15.17)$$

is also a Gaussian distribution.



2. If the prediction  $\mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t})$  is Gaussian and the sensor model  $\mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1})$  is linear Gaussian, then, after conditioning on the new evidence, the updated distribution

$$\mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) = \alpha \mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}) \quad (15.18)$$

is also a Gaussian distribution.

Thus, the FORWARD operator for Kalman filtering takes a Gaussian forward message  $\mathbf{f}_{1:t}$ , specified by a mean  $\boldsymbol{\mu}_t$  and covariance matrix  $\boldsymbol{\Sigma}_t$ , and produces a new multivariate Gaussian forward message  $\mathbf{f}_{1:t+1}$ , specified by a mean  $\boldsymbol{\mu}_{t+1}$  and covariance matrix  $\boldsymbol{\Sigma}_{t+1}$ . So, if we start with a Gaussian prior  $\mathbf{f}_{1:0} = \mathbf{P}(\mathbf{X}_0) = N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ , filtering with a linear Gaussian model produces a Gaussian state distribution for all time.



This seems to be a nice, elegant result, but why is it so important? The reason is that, except for a few special cases such as this, *filtering with continuous or hybrid (discrete and continuous) networks generates state distributions whose representation grows without bound over time*. This statement is not easy to prove in general, but Exercise 15.10 shows what happens for a simple example.

### 15.4.2 A simple one-dimensional example

We have said that the FORWARD operator for the Kalman filter maps a Gaussian into a new Gaussian. This translates into computing a new mean and covariance matrix from the previous mean and covariance matrix. Deriving the update rule in the general (multivariate) case requires rather a lot of linear algebra, so we will stick to a very simple univariate case for now; and later give the results for the general case. Even for the univariate case, the calculations are somewhat tedious, but we feel that they are worth seeing because the usefulness of the Kalman filter is tied so intimately to the mathematical properties of Gaussian distributions.

The temporal model we consider describes a **random walk** of a single continuous state variable  $X_t$  with a noisy observation  $Z_t$ . An example might be the “consumer confidence” index, which can be modeled as undergoing a random Gaussian-distributed change each month and is measured by a random consumer survey that also introduces Gaussian sampling noise.

The prior distribution is assumed to be Gaussian with variance  $\sigma_0^2$ :

$$P(x_0) = \alpha e^{-\frac{1}{2} \left( \frac{(x_0 - \mu_0)^2}{\sigma_0^2} \right)}.$$

(For simplicity, we use the same symbol  $\alpha$  for all normalizing constants in this section.) The transition model adds a Gaussian perturbation of constant variance  $\sigma_x^2$  to the current state:

$$P(x_{t+1} | x_t) = \alpha e^{-\frac{1}{2} \left( \frac{(x_{t+1} - x_t)^2}{\sigma_x^2} \right)}.$$

The sensor model assumes Gaussian noise with variance  $\sigma_z^2$ :

$$P(z_t | x_t) = \alpha e^{-\frac{1}{2} \left( \frac{(z_t - x_t)^2}{\sigma_z^2} \right)}.$$

Now, given the prior  $\mathbf{P}(X_0)$ , the one-step predicted distribution comes from Equation (15.17):

$$\begin{aligned} P(x_1) &= \int_{-\infty}^{\infty} P(x_1 | x_0) P(x_0) dx_0 = \alpha \int_{-\infty}^{\infty} e^{-\frac{1}{2} \left( \frac{(x_1 - x_0)^2}{\sigma_x^2} \right)} e^{-\frac{1}{2} \left( \frac{(x_0 - \mu_0)^2}{\sigma_0^2} \right)} dx_0 \\ &= \alpha \int_{-\infty}^{\infty} e^{-\frac{1}{2} \left( \frac{\sigma_0^2 (x_1 - x_0)^2 + \sigma_x^2 (x_0 - \mu_0)^2}{\sigma_0^2 \sigma_x^2} \right)} dx_0. \end{aligned}$$

This integral looks rather complicated. The key to progress is to notice that the exponent is the sum of two expressions that are *quadratic* in  $x_0$  and hence is itself a quadratic in  $x_0$ . A simple trick known as **completing the square** allows the rewriting of any quadratic  $ax_0^2 + bx_0 + c$  as the sum of a squared term  $a(x_0 - \frac{-b}{2a})^2$  and a residual term  $c - \frac{b^2}{4a}$  that is independent of  $x_0$ . The residual term can be taken outside the integral, giving us

$$P(x_1) = \alpha e^{-\frac{1}{2} \left( c - \frac{b^2}{4a} \right)} \int_{-\infty}^{\infty} e^{-\frac{1}{2} \left( a(x_0 - \frac{-b}{2a})^2 \right)} dx_0.$$

Now the integral is just the integral of a Gaussian over its full range, which is simply 1. Thus, we are left with only the residual term from the quadratic. Then, we notice that the residual term is a quadratic in  $x_1$ ; in fact, after simplification, we obtain

$$P(x_1) = \alpha e^{-\frac{1}{2} \left( \frac{(x_1 - \mu_0)^2}{\sigma_0^2 + \sigma_x^2} \right)}.$$

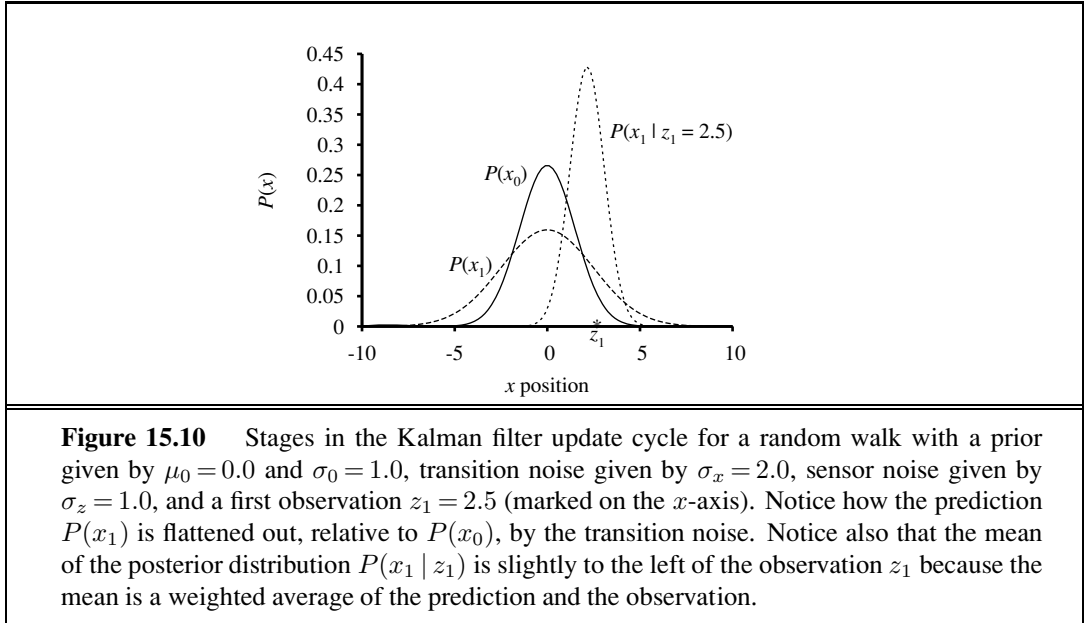
That is, the one-step predicted distribution is a Gaussian with the same mean  $\mu_0$  and a variance equal to the sum of the original variance  $\sigma_0^2$  and the transition variance  $\sigma_x^2$ .

To complete the update step, we need to condition on the observation at the first time step, namely,  $z_1$ . From Equation (15.18), this is given by

$$\begin{aligned} P(x_1 | z_1) &= \alpha P(z_1 | x_1) P(x_1) \\ &= \alpha e^{-\frac{1}{2} \left( \frac{(z_1 - x_1)^2}{\sigma_z^2} \right)} e^{-\frac{1}{2} \left( \frac{(x_1 - \mu_0)^2}{\sigma_0^2 + \sigma_x^2} \right)}. \end{aligned}$$

Once again, we combine the exponents and complete the square (Exercise 15.11), obtaining

$$P(x_1 | z_1) = \alpha e^{-\frac{1}{2} \left( \frac{(x_1 - \frac{(\sigma_0^2 + \sigma_x^2)z_1 + \sigma_z^2 \mu_0}{\sigma_0^2 + \sigma_x^2 + \sigma_z^2})^2}{(\frac{\sigma_0^2 + \sigma_x^2}{\sigma_0^2 + \sigma_x^2 + \sigma_z^2}) \sigma_z^2 / (\frac{\sigma_0^2 + \sigma_x^2}{\sigma_0^2 + \sigma_x^2 + \sigma_z^2} + \sigma_z^2)} \right)}. \quad (15.19)$$



Thus, after one update cycle, we have a new Gaussian distribution for the state variable.

From the Gaussian formula in Equation (15.19), we see that the new mean and standard deviation can be calculated from the old mean and standard deviation as follows:

$$\mu_{t+1} = \frac{(\sigma_t^2 + \sigma_x^2)z_{t+1} + \sigma_z^2\mu_t}{\sigma_t^2 + \sigma_x^2 + \sigma_z^2} \quad \text{and} \quad \sigma_{t+1}^2 = \frac{(\sigma_t^2 + \sigma_x^2)\sigma_z^2}{\sigma_t^2 + \sigma_x^2 + \sigma_z^2}. \quad (15.20)$$

Figure 15.10 shows one update cycle for particular values of the transition and sensor models.

Equation (15.20) plays exactly the same role as the general filtering equation (15.5) or the HMM filtering equation (15.12). Because of the special nature of Gaussian distributions, however, the equations have some interesting additional properties. First, we can interpret the calculation for the new mean  $\mu_{t+1}$  as simply a *weighted mean* of the new observation  $z_{t+1}$  and the old mean  $\mu_t$ . If the observation is unreliable, then  $\sigma_z^2$  is large and we pay more attention to the old mean; if the old mean is unreliable ( $\sigma_t^2$  is large) or the process is highly unpredictable ( $\sigma_x^2$  is large), then we pay more attention to the observation. Second, notice that the update for the variance  $\sigma_{t+1}^2$  is *independent of the observation*. We can therefore compute in advance what the sequence of variance values will be. Third, the sequence of variance values converges quickly to a fixed value that depends only on  $\sigma_x^2$  and  $\sigma_z^2$ , thereby substantially simplifying the subsequent calculations. (See Exercise 15.12.)

### 15.4.3 The general case

The preceding derivation illustrates the key property of Gaussian distributions that allows Kalman filtering to work: the fact that the exponent is a quadratic form. This is true not just for the univariate case; the full multivariate Gaussian distribution has the form

$$N(\boldsymbol{\mu}, \boldsymbol{\Sigma})(\mathbf{x}) = \alpha e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}.$$

Multiplying out the terms in the exponent makes it clear that the exponent is also a quadratic function of the values  $x_i$  in  $\mathbf{x}$ . As in the univariate case, the filtering update preserves the Gaussian nature of the state distribution.

Let us first define the general temporal model used with Kalman filtering. Both the transition model and the sensor model allow for a *linear* transformation with additive Gaussian noise. Thus, we have

$$\begin{aligned} P(\mathbf{x}_{t+1} | \mathbf{x}_t) &= N(\mathbf{F}\mathbf{x}_t, \Sigma_x)(\mathbf{x}_{t+1}) \\ P(\mathbf{z}_t | \mathbf{x}_t) &= N(\mathbf{H}\mathbf{x}_t, \Sigma_z)(\mathbf{z}_t), \end{aligned} \quad (15.21)$$

where  $\mathbf{F}$  and  $\Sigma_x$  are matrices describing the linear transition model and transition noise covariance, and  $\mathbf{H}$  and  $\Sigma_z$  are the corresponding matrices for the sensor model. Now the update equations for the mean and covariance, in their full, hairy horribleness, are

$$\begin{aligned} \mu_{t+1} &= \mathbf{F}\mu_t + \mathbf{K}_{t+1}(\mathbf{z}_{t+1} - \mathbf{H}\mathbf{F}\mu_t) \\ \Sigma_{t+1} &= (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{H})(\mathbf{F}\Sigma_t\mathbf{F}^\top + \Sigma_x), \end{aligned} \quad (15.22)$$

where  $\mathbf{K}_{t+1} = (\mathbf{F}\Sigma_t\mathbf{F}^\top + \Sigma_x)\mathbf{H}^\top(\mathbf{H}(\mathbf{F}\Sigma_t\mathbf{F}^\top + \Sigma_x)\mathbf{H}^\top + \Sigma_z)^{-1}$  is called the **Kalman gain matrix**. Believe it or not, these equations make some intuitive sense. For example, consider the update for the mean state estimate  $\mu$ . The term  $\mathbf{F}\mu_t$  is the *predicted* state at  $t + 1$ , so  $\mathbf{H}\mathbf{F}\mu_t$  is the *predicted* observation. Therefore, the term  $\mathbf{z}_{t+1} - \mathbf{H}\mathbf{F}\mu_t$  represents the error in the predicted observation. This is multiplied by  $\mathbf{K}_{t+1}$  to correct the predicted state; hence,  $\mathbf{K}_{t+1}$  is a measure of *how seriously to take the new observation* relative to the prediction. As in Equation (15.20), we also have the property that the variance update is independent of the observations. The sequence of values for  $\Sigma_t$  and  $\mathbf{K}_t$  can therefore be computed offline, and the actual calculations required during online tracking are quite modest.

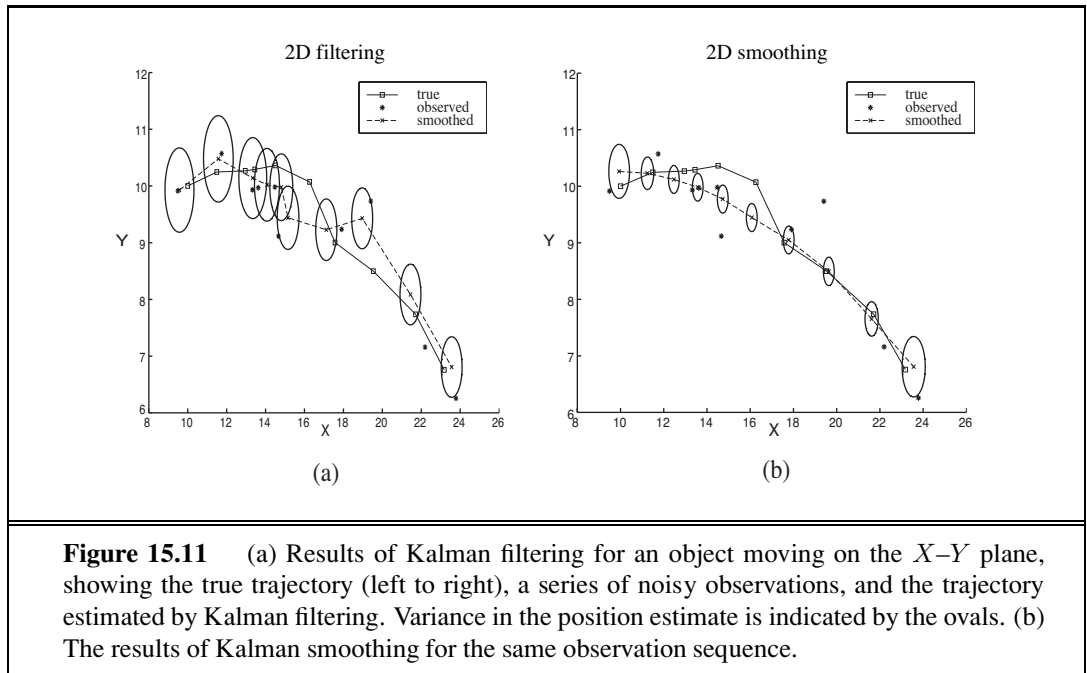
To illustrate these equations at work, we have applied them to the problem of tracking an object moving on the  $X$ - $Y$  plane. The state variables are  $\mathbf{X} = (X, Y, \dot{X}, \dot{Y})^\top$ , so  $\mathbf{F}$ ,  $\Sigma_x$ ,  $\mathbf{H}$ , and  $\Sigma_z$  are  $4 \times 4$  matrices. Figure 15.11(a) shows the true trajectory, a series of noisy observations, and the trajectory estimated by Kalman filtering, along with the covariances indicated by the one-standard-deviation contours. The filtering process does a good job of tracking the actual motion, and, as expected, the variance quickly reaches a fixed point.

We can also derive equations for *smoothing* as well as filtering with linear Gaussian models. The smoothing results are shown in Figure 15.11(b). Notice how the variance in the position estimate is sharply reduced, except at the ends of the trajectory (why?), and that the estimated trajectory is much smoother.

#### 15.4.4 Applicability of Kalman filtering

The Kalman filter and its elaborations are used in a vast array of applications. The “classical” application is in radar tracking of aircraft and missiles. Related applications include acoustic tracking of submarines and ground vehicles and visual tracking of vehicles and people. In a slightly more esoteric vein, Kalman filters are used to reconstruct particle trajectories from bubble-chamber photographs and ocean currents from satellite surface measurements. The range of application is much larger than just the tracking of motion: any system characterized by continuous state variables and noisy measurements will do. Such systems include pulp mills, chemical plants, nuclear reactors, plant ecosystems, and national economies.





**Figure 15.11** (a) Results of Kalman filtering for an object moving on the  $X$ - $Y$  plane, showing the true trajectory (left to right), a series of noisy observations, and the trajectory estimated by Kalman filtering. Variance in the position estimate is indicated by the ovals. (b) The results of Kalman smoothing for the same observation sequence.

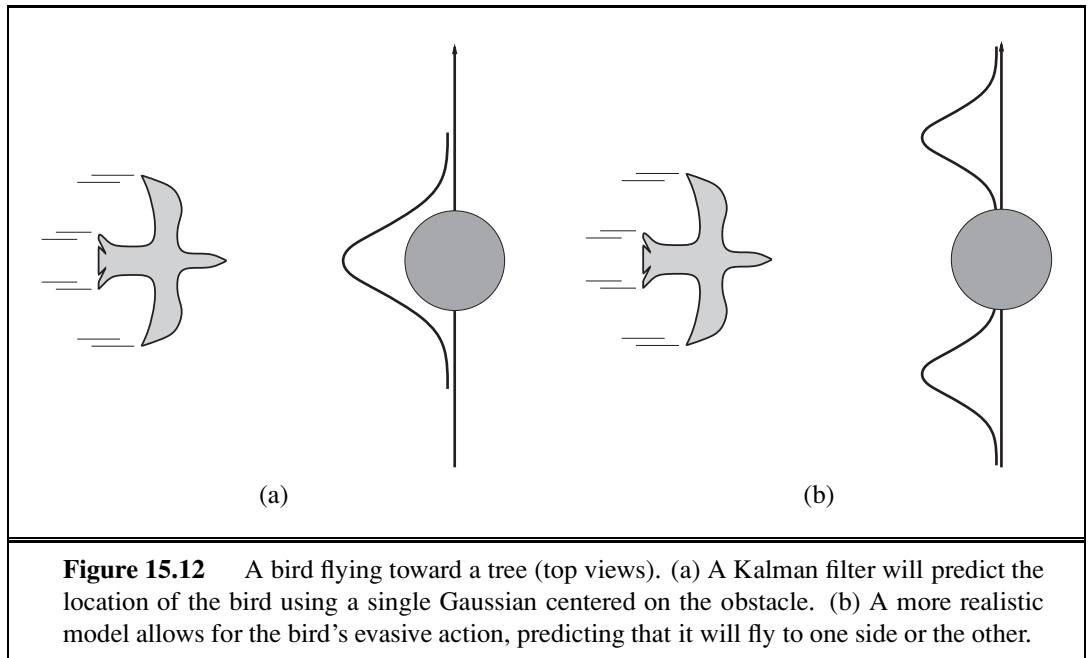
The fact that Kalman filtering can be applied to a system does not mean that the results will be valid or useful. The assumptions made—a linear Gaussian transition and sensor models—are very strong. The **extended Kalman filter (EKF)** attempts to overcome nonlinearities in the system being modeled. A system is **nonlinear** if the transition model cannot be described as a matrix multiplication of the state vector, as in Equation (15.21). The EKF works by modeling the system as *locally* linear in  $\mathbf{x}_t$  in the region of  $\mathbf{x}_t = \boldsymbol{\mu}_t$ , the mean of the current state distribution. This works well for smooth, well-behaved systems and allows the tracker to maintain and update a Gaussian state distribution that is a reasonable approximation to the true posterior. A detailed example is given in Chapter 25.

What does it mean for a system to be “unsmooth” or “poorly behaved”? Technically, it means that there is significant nonlinearity in system response within the region that is “close” (according to the covariance  $\boldsymbol{\Sigma}_t$ ) to the current mean  $\boldsymbol{\mu}_t$ . To understand this idea in nontechnical terms, consider the example of trying to track a bird as it flies through the jungle. The bird appears to be heading at high speed straight for a tree trunk. The Kalman filter, whether regular or extended, can make only a Gaussian prediction of the location of the bird, and the mean of this Gaussian will be centered on the trunk, as shown in Figure 15.12(a). A reasonable model of the bird, on the other hand, would predict evasive action to one side or the other, as shown in Figure 15.12(b). Such a model is highly nonlinear, because the bird’s decision varies sharply depending on its precise location relative to the trunk.

To handle examples like these, we clearly need a more expressive language for representing the behavior of the system being modeled. Within the control theory community, for which problems such as evasive maneuvering by aircraft raise the same kinds of difficulties, the standard solution is the **switching Kalman filter**. In this approach, multiple Kalman fil-

EXTENDED KALMAN  
FILTER (EKF)  
NONLINEAR

SWITCHING KALMAN  
FILTER



**Figure 15.12** A bird flying toward a tree (top views). (a) A Kalman filter will predict the location of the bird using a single Gaussian centered on the obstacle. (b) A more realistic model allows for the bird's evasive action, predicting that it will fly to one side or the other.

ters run in parallel, each using a different model of the system—for example, one for straight flight, one for sharp left turns, and one for sharp right turns. A weighted sum of predictions is used, where the weight depends on how well each filter fits the current data. We will see in the next section that this is simply a special case of the general dynamic Bayesian network model, obtained by adding a discrete “maneuver” state variable to the network shown in Figure 15.9. Switching Kalman filters are discussed further in Exercise 15.10.

## 15.5 DYNAMIC BAYESIAN NETWORKS

### DYNAMIC BAYESIAN NETWORK

A **dynamic Bayesian network**, or **DBN**, is a Bayesian network that represents a temporal probability model of the kind described in Section 15.1. We have already seen examples of DBNs: the umbrella network in Figure 15.2 and the Kalman filter network in Figure 15.9. In general, each slice of a DBN can have any number of state variables  $\mathbf{X}_t$  and evidence variables  $\mathbf{E}_t$ . For simplicity, we assume that the variables and their links are exactly replicated from slice to slice and that the DBN represents a first-order Markov process, so that each variable can have parents only in its own slice or the immediately preceding slice.

It should be clear that every hidden Markov model can be represented as a DBN with a single state variable and a single evidence variable. It is also the case that every discrete-variable DBN can be represented as an HMM; as explained in Section 15.3, we can combine all the state variables in the DBN into a single state variable whose values are all possible tuples of values of the individual state variables. Now, if every HMM is a DBN and every DBN can be translated into an HMM, what's the difference? The difference is that, *by de-*



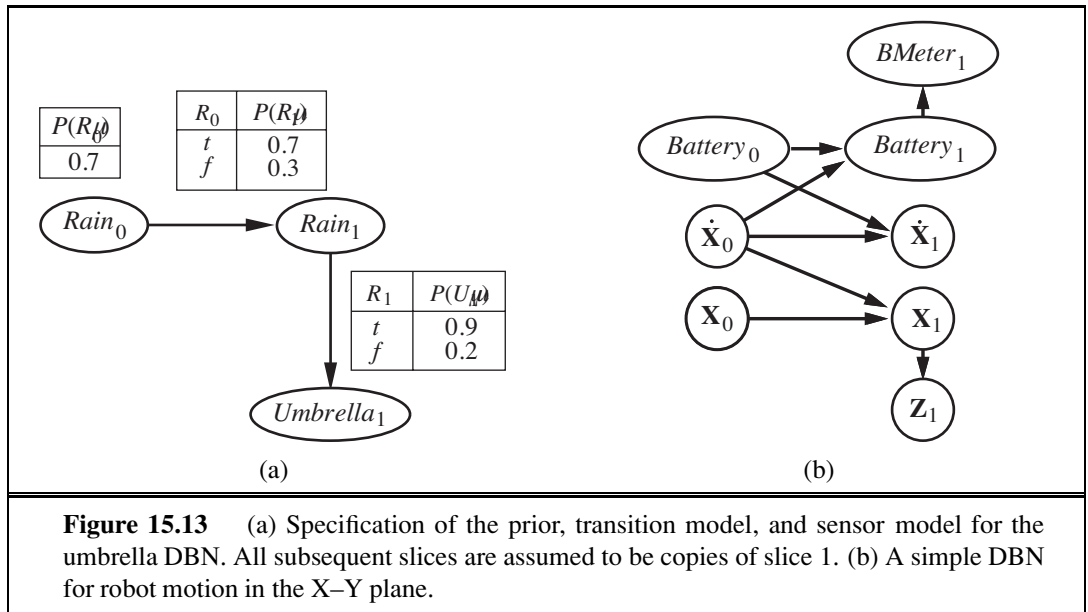
*composing the state of a complex system into its constituent variables, the can take advantage of sparseness in the temporal probability model.* Suppose, for example, that a DBN has 20 Boolean state variables, each of which has three parents in the preceding slice. Then the DBN transition model has  $20 \times 2^3 = 160$  probabilities, whereas the corresponding HMM has  $2^{20}$  states and therefore  $2^{40}$ , or roughly a trillion, probabilities in the transition matrix. This is bad for at least three reasons: first, the HMM itself requires much more space; second, the huge transition matrix makes HMM inference much more expensive; and third, the problem of learning such a huge number of parameters makes the pure HMM model unsuitable for large problems. The relationship between DBNs and HMMs is roughly analogous to the relationship between ordinary Bayesian networks and full tabulated joint distributions.

We have already explained that every Kalman filter model can be represented in a DBN with continuous variables and linear Gaussian conditional distributions (Figure 15.9). It should be clear from the discussion at the end of the preceding section that *not* every DBN can be represented by a Kalman filter model. In a Kalman filter, the current state distribution is always a single multivariate Gaussian distribution—that is, a single “bump” in a particular location. DBNs, on the other hand, can model arbitrary distributions. For many real-world applications, this flexibility is essential. Consider, for example, the current location of my keys. They might be in my pocket, on the bedside table, on the kitchen counter, dangling from the front door, or locked in the car. A single Gaussian bump that included all these places would have to allocate significant probability to the keys being in mid-air in the front hall. Aspects of the real world such as purposive agents, obstacles, and pockets introduce “nonlinearities” that require combinations of discrete and continuous variables in order to get reasonable models.

### 15.5.1 Constructing DBNs

To construct a DBN, one must specify three kinds of information: the prior distribution over the state variables,  $\mathbf{P}(\mathbf{X}_0)$ ; the transition model  $\mathbf{P}(\mathbf{X}_{t+1} \mid \mathbf{X}_t)$ ; and the sensor model  $\mathbf{P}(\mathbf{E}_t \mid \mathbf{X}_t)$ . To specify the transition and sensor models, one must also specify the topology of the connections between successive slices and between the state and evidence variables. Because the transition and sensor models are assumed to be stationary—the same for all  $t$ —it is most convenient simply to specify them for the first slice. For example, the complete DBN specification for the umbrella world is given by the three-node network shown in Figure 15.13(a). From this specification, the complete DBN with an unbounded number of time slices can be constructed as needed by copying the first slice.

Let us now consider a more interesting example: monitoring a battery-powered robot moving in the  $X$ - $Y$  plane, as introduced at the end of Section 15.1. First, we need state variables, which will include both  $\mathbf{X}_t = (X_t, Y_t)$  for position and  $\dot{\mathbf{X}}_t = (\dot{X}_t, \dot{Y}_t)$  for velocity. We assume some method of measuring position—perhaps a fixed camera or onboard GPS (Global Positioning System)—yielding measurements  $\mathbf{Z}_t$ . The position at the next time step depends on the current position and velocity, as in the standard Kalman filter model. The velocity at the next step depends on the current velocity and the state of the battery. We add  $Battery_t$  to represent the actual battery charge level, which has as parents the previous



battery level and the velocity, and we add  $BMeter_t$ , which measures the battery charge level. This gives us the basic model shown in Figure 15.13(b).

It is worth looking in more depth at the nature of the sensor model for  $BMeter_t$ . Let us suppose, for simplicity, that both  $Battery_t$  and  $BMeter_t$  can take on discrete values 0 through 5. If the meter is always accurate, then the CPT  $P(BMeter_t | Battery_t)$  should have probabilities of 1.0 “along the diagonal” and probabilities of 0.0 elsewhere. In reality, noise always creeps into measurements. For continuous measurements, a Gaussian distribution with a small variance might be used.<sup>5</sup> For our discrete variables, we can approximate a Gaussian using a distribution in which the probability of error drops off in the appropriate way, so that the probability of a large error is very small. We use the term **Gaussian error model** to cover both the continuous and discrete versions.

Anyone with hands-on experience of robotics, computerized process control, or other forms of automatic sensing will readily testify to the fact that small amounts of measurement noise are often the least of one’s problems. Real sensors *fail*. When a sensor fails, it does not necessarily send a signal saying, “Oh, by the way, the data I’m about to send you is a load of nonsense.” Instead, it simply sends the nonsense. The simplest kind of failure is called a **transient failure**, where the sensor occasionally decides to send some nonsense. For example, the battery level sensor might have a habit of sending a zero when someone bumps the robot, even if the battery is fully charged.

Let’s see what happens when a transient failure occurs with a Gaussian error model that doesn’t accommodate such failures. Suppose, for example, that the robot is sitting quietly and observes 20 consecutive battery readings of 5. Then the battery meter has a temporary seizure

<sup>5</sup> Strictly speaking, a Gaussian distribution is problematic because it assigns nonzero probability to large negative charge levels. The **beta distribution** is sometimes a better choice for a variable whose range is restricted.

and the next reading is  $BMeter_{21} = 0$ . What will the simple Gaussian error model lead us to believe about  $Battery_{21}$ ? According to Bayes' rule, the answer depends on both the sensor model  $\mathbf{P}(BMeter_{21} = 0 \mid Battery_{21})$  and the prediction  $\mathbf{P}(Battery_{21} \mid BMeter_{1:20})$ . If the probability of a large sensor error is significantly less likely than the probability of a transition to  $Battery_{21} = 0$ , even if the latter is very unlikely, then the posterior distribution will assign a high probability to the battery's being empty. A second reading of 0 at  $t = 22$  will make this conclusion almost certain. If the transient failure then disappears and the reading returns to 5 from  $t = 23$  onwards, the estimate for the battery level will quickly return to 5, as if by magic. This course of events is illustrated in the upper curve of Figure 15.14(a), which shows the expected value of  $Battery_t$  over time, using a discrete Gaussian error model.

Despite the recovery, there is a time ( $t = 22$ ) when the robot is convinced that its battery is empty; presumably, then, it should send out a mayday signal and shut down. Alas, its oversimplified sensor model has led it astray. How can this be fixed? Consider a familiar example from everyday human driving: on sharp curves or steep hills, one's "fuel tank empty" warning light sometimes turns on. Rather than looking for the emergency phone, one simply recalls that the fuel gauge sometimes gives a very large error when the fuel is sloshing around in the tank. The moral of the story is the following: *for the system to handle sensor failure properly, the sensor model must include the possibility of failure.*

The simplest kind of failure model for a sensor allows a certain probability that the sensor will return some completely incorrect value, regardless of the true state of the world. For example, if the battery meter fails by returning 0, we might say that

$$P(BMeter_t = 0 \mid Battery_t = 5) = 0.03 ,$$

which is presumably much larger than the probability assigned by the simple Gaussian error model. Let's call this the **transient failure model**. How does it help when we are faced with a reading of 0? Provided that the *predicted* probability of an empty battery, according to the readings so far, is much less than 0.03, then the best explanation of the observation  $BMeter_{21} = 0$  is that the sensor has temporarily failed. Intuitively, we can think of the belief about the battery level as having a certain amount of "inertia" that helps to overcome temporary blips in the meter reading. The upper curve in Figure 15.14(b) shows that the transient failure model can handle transient failures without a catastrophic change in beliefs.

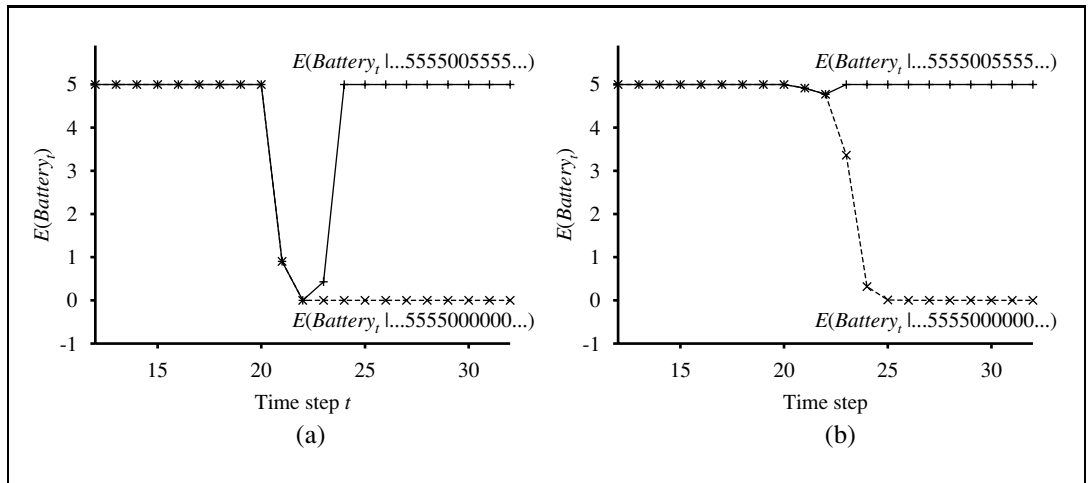
So much for temporary blips. What about a persistent sensor failure? Sadly, failures of this kind are all too common. If the sensor returns 20 readings of 5 followed by 20 readings of 0, then the transient sensor failure model described in the preceding paragraph will result in the robot gradually coming to believe that its battery is empty when in fact it may be that the meter has failed. The lower curve in Figure 15.14(b) shows the belief "trajectory" for this case. By  $t = 25$ —five readings of 0—the robot is convinced that its battery is empty. Obviously, we would prefer the robot to believe that its battery meter is broken—if indeed this is the more likely event.

Unsurprisingly, to handle persistent failure, we need a **persistent failure model** that describes how the sensor behaves under normal conditions and after failure. To do this, we need to augment the state of the system with an additional variable, say,  $BMBroken$ , that describes the status of the battery meter. The persistence of failure must be modeled by an

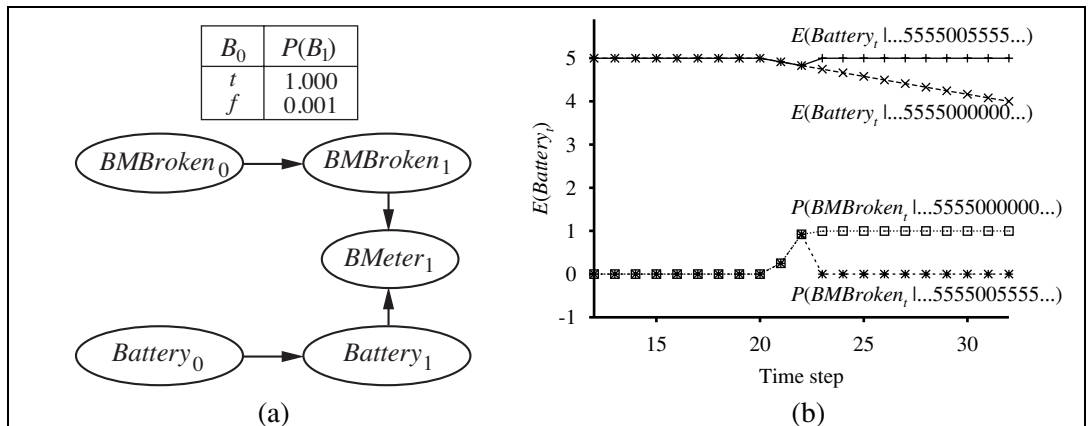


TRANSIENT FAILURE  
MODEL

PERSISTENT  
FAILURE MODEL



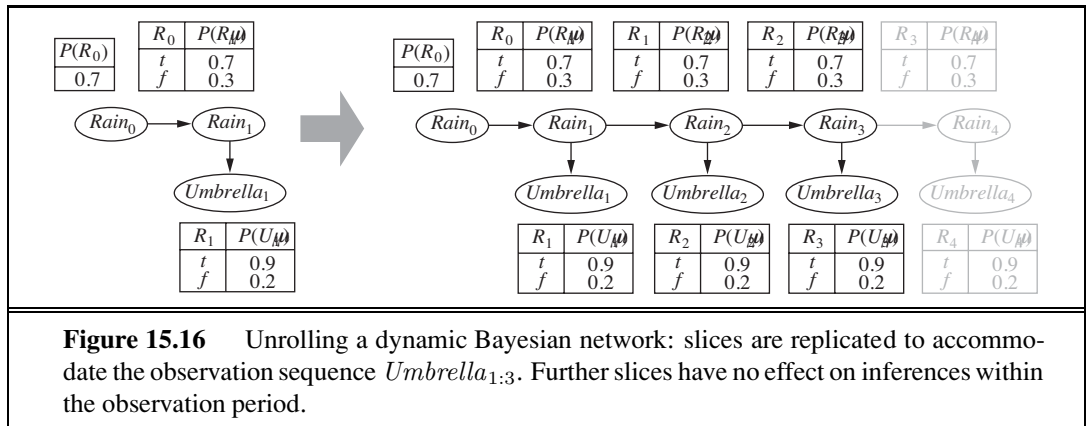
**Figure 15.14** (a) Upper curve: trajectory of the expected value of  $Battery_t$  for an observation sequence consisting of all 5s except for 0s at  $t = 21$  and  $t = 22$ , using a simple Gaussian error model. Lower curve: trajectory when the observation remains at 0 from  $t = 21$  onwards. (b) The same experiment run with the transient failure model. Notice that the transient failure is handled well, but the persistent failure results in excessive pessimism about the battery charge.



**Figure 15.15** (a) A DBN fragment showing the sensor status variable required for modeling persistent failure of the battery sensor. (b) Upper curves: trajectories of the expected value of  $Battery_t$  for the "transient failure" and "permanent failure" observation sequences. Lower curves: probability trajectories for  $BMBroken$  given the two observation sequences.

#### PERSISTENCE ARC

arc linking  $BMBroken_0$  to  $BMBroken_1$ . This **persistence arc** has a CPT that gives a small probability of failure in any given time step, say, 0.001, but specifies that the sensor stays broken once it breaks. When the sensor is OK, the sensor model for  $BMeter$  is identical to the transient failure model; when the sensor is broken, it says  $BMeter$  is always 0, regardless of the actual battery charge.



The persistent failure model for the battery sensor is shown in Figure 15.15(a). Its performance on the two data sequences (temporary blip and persistent failure) is shown in Figure 15.15(b). There are several things to notice about these curves. First, in the case of the temporary blip, the probability that the sensor is broken rises significantly after the second 0 reading, but immediately drops back to zero once a 5 is observed. Second, in the case of persistent failure, the probability that the sensor is broken rises quickly to almost 1 and stays there. Finally, once the sensor is known to be broken, the robot can only assume that its battery discharges at the “normal” rate, as shown by the gradually descending level of  $E(Battery_t | \dots)$ .

So far, we have merely scratched the surface of the problem of representing complex processes. The variety of transition models is huge, encompassing topics as disparate as modeling the human endocrine system and modeling multiple vehicles driving on a freeway. Sensor modeling is also a vast subfield in itself, but even subtle phenomena, such as sensor drift, sudden decalibration, and the effects of exogenous conditions (such as weather) on sensor readings, can be handled by explicit representation within dynamic Bayesian networks.

### 15.5.2 Exact inference in DBNs

Having sketched some ideas for representing complex processes as DBNs, we now turn to the question of inference. In a sense, this question has already been answered: dynamic Bayesian networks *are* Bayesian networks, and we already have algorithms for inference in Bayesian networks. Given a sequence of observations, one can construct the full Bayesian network representation of a DBN by replicating slices until the network is large enough to accommodate the observations, as in Figure 15.16. This technique, mentioned in Chapter 14 in the context of relational probability models, is called **unrolling**. (Technically, the DBN is equivalent to the semi-infinite network obtained by unrolling forever. Slices added beyond the last observation have no effect on inferences within the observation period and can be omitted.) Once the DBN is unrolled, one can use any of the inference algorithms—variable elimination, clustering methods, and so on—described in Chapter 14.

Unfortunately, a naive application of unrolling would not be particularly efficient. If we want to perform filtering or smoothing with a long sequence of observations  $\mathbf{e}_{1:t}$ , the

unrolled network would require  $O(t)$  space and would thus grow without bound as more observations were added. Moreover, if we simply run the inference algorithm anew each time an observation is added, the inference time per update will also increase as  $O(t)$ .

Looking back to Section 15.2.1, we see that constant time and space per filtering update can be achieved if the computation can be done recursively. Essentially, the filtering update in Equation (15.5) works by *summing out* the state variables of the previous time step to get the distribution for the new time step. Summing out variables is exactly what the **variable elimination** (Figure 14.11) algorithm does, and it turns out that running variable elimination with the variables in temporal order exactly mimics the operation of the recursive filtering update in Equation (15.5). The modified algorithm keeps at most two slices in memory at any one time: starting with slice 0, we add slice 1, then sum out slice 0, then add slice 2, then sum out slice 1, and so on. In this way, we can achieve constant space and time per filtering update. (The same performance can be achieved by suitable modifications to the clustering algorithm.) Exercise 15.17 asks you to verify this fact for the umbrella network.

So much for the good news; now for the bad news: It turns out that the “constant” for the per-update time and space complexity is, in almost all cases, exponential in the number of state variables. What happens is that, as the variable elimination proceeds, the factors grow to include all the state variables (or, more precisely, all those state variables that have parents in the previous time slice). The maximum factor size is  $O(d^{n+k})$  and the total update cost per step is  $O(nd^{n+k})$ , where  $d$  is the domain size of the variables and  $k$  is the maximum number of parents of any state variable.

Of course, this is much less than the cost of HMM updating, which is  $O(d^{2n})$ , but it is still infeasible for large numbers of variables. This grim fact is somewhat hard to accept. What it means is that *even though we can use DBNs to represent very complex temporal processes with many sparsely connected variables, we cannot reason efficiently and exactly about those processes*. The DBN model itself, which represents the prior joint distribution over all the variables, is factorable into its constituent CPTs, but the posterior joint distribution conditioned on an observation sequence—that is, the forward message—is generally *not* factorable. So far, no one has found a way around this problem, despite the fact that many important areas of science and engineering would benefit enormously from its solution. Thus, we must fall back on approximate methods.



### 15.5.3 Approximate inference in DBNs

Section 14.5 described two approximation algorithms: likelihood weighting (Figure 14.15) and Markov chain Monte Carlo (MCMC, Figure 14.16). Of the two, the former is most easily adapted to the DBN context. (An MCMC filtering algorithm is described briefly in the notes at the end of the chapter.) We will see, however, that several improvements are required over the standard likelihood weighting algorithm before a practical method emerges.

Recall that likelihood weighting works by sampling the nonevidence nodes of the network in topological order, weighting each sample by the likelihood it accords to the observed evidence variables. As with the exact algorithms, we could apply likelihood weighting directly to an unrolled DBN, but this would suffer from the same problems of increasing time





and space requirements per update as the observation sequence grows. The problem is that the standard algorithm runs each sample in turn, all the way through the network. Instead, we can simply run all  $N$  samples together through the DBN, one slice at a time. The modified algorithm fits the general pattern of filtering algorithms, with the set of  $N$  samples as the forward message. The first key innovation, then, is to *use the samples themselves as an approximate representation of the current state distribution*. This meets the requirement of a “constant” time per update, although the constant depends on the number of samples required to maintain an accurate approximation. There is also no need to unroll the DBN, because we need to have in memory only the current slice and the next slice.

In our discussion of likelihood weighting in Chapter 14, we pointed out that the algorithm’s accuracy suffers if the evidence variables are “downstream” from the variables being sampled, because in that case the samples are generated without any influence from the evidence. Looking at the typical structure of a DBN—say, the umbrella DBN in Figure 15.16—we see that indeed the early state variables will be sampled without the benefit of the later evidence. In fact, looking more carefully, we see that *none* of the state variables has *any* evidence variables among its ancestors! Hence, although the weight of each sample will depend on the evidence, the actual set of samples generated will be *completely independent* of the evidence. For example, even if the boss brings in the umbrella every day, the sampling process could still hallucinate endless days of sunshine. What this means in practice is that the fraction of samples that remain reasonably close to the actual series of events (and therefore have nonnegligible weights) drops exponentially with  $t$ , the length of the observation sequence. In other words, to maintain a given level of accuracy, we need to increase the number of samples exponentially with  $t$ . Given that a filtering algorithm that works in real time can use only a fixed number of samples, what happens in practice is that the error blows up after a very small number of update steps.



Clearly, we need a better solution. The second key innovation is to *focus the set of samples on the high-probability regions of the state space*. This can be done by throwing away samples that have very low weight, according to the observations, while replicating those that have high weight. In that way, the population of samples will stay reasonably close to reality. If we think of samples as a resource for modeling the posterior distribution, then it makes sense to use more samples in regions of the state space where the posterior is higher.

#### PARTICLE FILTERING

A family of algorithms called **particle filtering** is designed to do just that. Particle filtering works as follows: First, a population of  $N$  initial-state samples is created by sampling from the prior distribution  $\mathbf{P}(\mathbf{X}_0)$ . Then the update cycle is repeated for each time step:

1. Each sample is propagated forward by sampling the next state value  $\mathbf{x}_{t+1}$  given the current value  $\mathbf{x}_t$  for the sample, based on the transition model  $\mathbf{P}(\mathbf{X}_{t+1} \mid \mathbf{x}_t)$ .
2. Each sample is weighted by the likelihood it assigns to the new evidence,  $P(\mathbf{e}_{t+1} \mid \mathbf{x}_{t+1})$ .
3. The population is *resampled* to generate a new population of  $N$  samples. Each new sample is selected from the current population; the probability that a particular sample is selected is proportional to its weight. The new samples are unweighted.

The algorithm is shown in detail in Figure 15.17, and its operation for the umbrella DBN is illustrated in Figure 15.18.

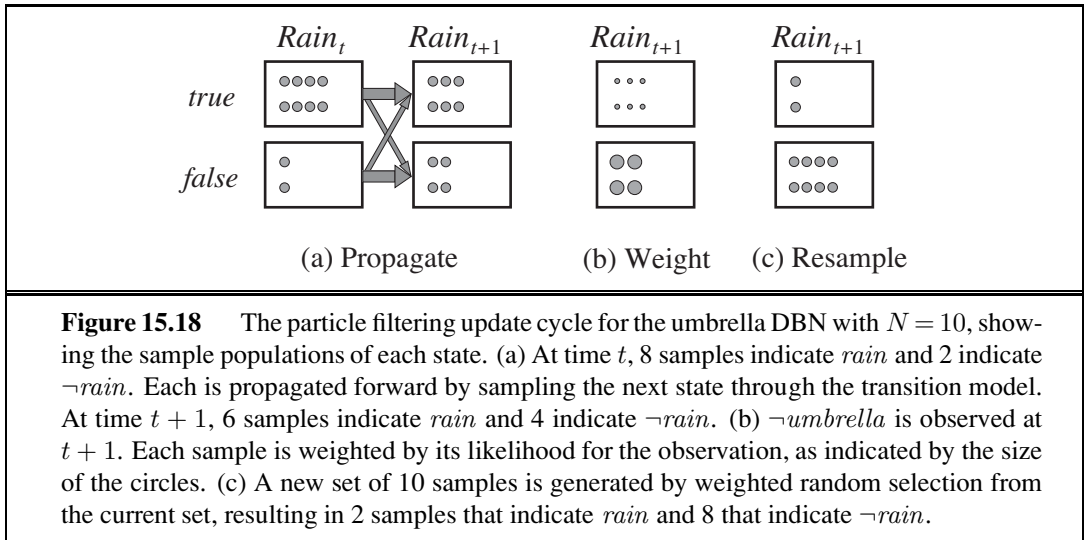
```

function PARTICLE-FILTERING( $\mathbf{e}, N, dbn$ ) returns a set of samples for the next time step
  inputs:  $\mathbf{e}$ , the new incoming evidence
            $N$ , the number of samples to be maintained
            $dbn$ , a DBN with prior  $\mathbf{P}(\mathbf{X}_0)$ , transition model  $\mathbf{P}(\mathbf{X}_1|\mathbf{X}_0)$ , sensor model  $\mathbf{P}(\mathbf{E}_1|\mathbf{X}_1)$ 
  persistent:  $S$ , a vector of samples of size  $N$ , initially generated from  $\mathbf{P}(\mathbf{X}_0)$ 
  local variables:  $W$ , a vector of weights of size  $N$ 

  for  $i = 1$  to  $N$  do
     $S[i] \leftarrow$  sample from  $\mathbf{P}(\mathbf{X}_1 | \mathbf{X}_0 = S[i])$  /* step 1 */
     $W[i] \leftarrow \mathbf{P}(\mathbf{e} | \mathbf{X}_1 = S[i])$  /* step 2 */
   $S \leftarrow$  WEIGHTED-SAMPLE-WITH-REPLACEMENT( $N, S, W$ ) /* step 3 */
  return  $S$ 

```

**Figure 15.17** The particle filtering algorithm implemented as a recursive update operation with state (the set of samples). Each of the sampling operations involves sampling the relevant slice variables in topological order, much as in PRIOR-SAMPLE. The WEIGHTED-SAMPLE-WITH-REPLACEMENT operation can be implemented to run in  $O(N)$  expected time. The step numbers refer to the description in the text.



**Figure 15.18** The particle filtering update cycle for the umbrella DBN with  $N = 10$ , showing the sample populations of each state. (a) At time  $t$ , 8 samples indicate *rain* and 2 indicate  $\neg$ *rain*. Each is propagated forward by sampling the next state through the transition model. At time  $t + 1$ , 6 samples indicate *rain* and 4 indicate  $\neg$ *rain*. (b)  $\neg$ *umbrella* is observed at  $t + 1$ . Each sample is weighted by its likelihood for the observation, as indicated by the size of the circles. (c) A new set of 10 samples is generated by weighted random selection from the current set, resulting in 2 samples that indicate *rain* and 8 that indicate  $\neg$ *rain*.

We can show that this algorithm is consistent—gives the correct probabilities as  $N$  tends to infinity—by considering what happens during one update cycle. We assume that the sample population starts with a correct representation of the forward message  $\mathbf{f}_{1:t} = \mathbf{P}(\mathbf{X}_t | \mathbf{e}_{1:t})$  at time  $t$ . Writing  $N(\mathbf{x}_t | \mathbf{e}_{1:t})$  for the number of samples occupying state  $\mathbf{x}_t$  after observations  $\mathbf{e}_{1:t}$  have been processed, we therefore have

$$N(\mathbf{x}_t | \mathbf{e}_{1:t})/N = P(\mathbf{x}_t | \mathbf{e}_{1:t}) \quad (15.23)$$

for large  $N$ . Now we propagate each sample forward by sampling the state variables at  $t + 1$ , given the values for the sample at  $t$ . The number of samples reaching state  $\mathbf{x}_{t+1}$  from each

$\mathbf{x}_t$  is the transition probability times the population of  $\mathbf{x}_t$ ; hence, the total number of samples reaching  $\mathbf{x}_{t+1}$  is

$$N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t}) = \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1} | \mathbf{x}_t) N(\mathbf{x}_t | \mathbf{e}_{1:t}) .$$

Now we weight each sample by its likelihood for the evidence at  $t + 1$ . A sample in state  $\mathbf{x}_{t+1}$  receives weight  $P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1})$ . The total weight of the samples in  $\mathbf{x}_{t+1}$  after seeing  $\mathbf{e}_{t+1}$  is therefore

$$W(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) = P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t}) .$$

Now for the resampling step. Since each sample is replicated with probability proportional to its weight, the number of samples in state  $\mathbf{x}_{t+1}$  after resampling is proportional to the total weight in  $\mathbf{x}_{t+1}$  before resampling:

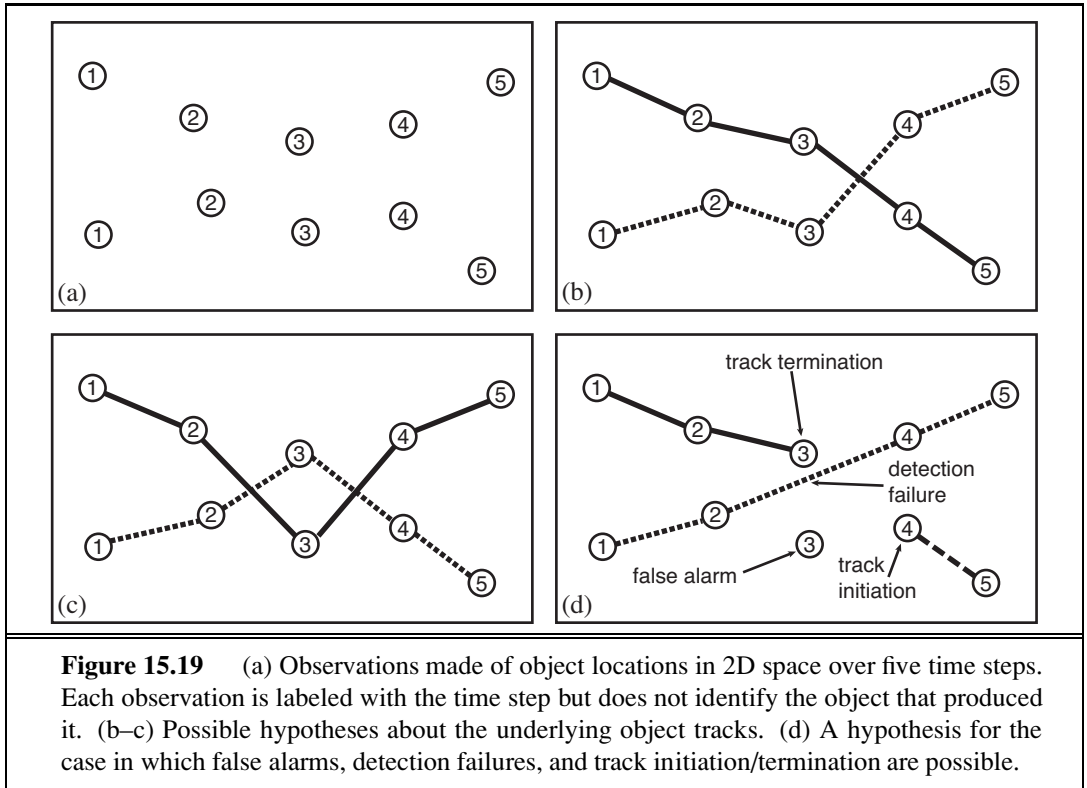
$$\begin{aligned} N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) / N &= \alpha W(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) \\ &= \alpha P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t}) \\ &= \alpha P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1} | \mathbf{x}_t) N(\mathbf{x}_t | \mathbf{e}_{1:t}) \\ &= \alpha N P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t}) \quad (\text{by 15.23}) \\ &= \alpha' P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t}) \\ &= P(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) \quad (\text{by 15.5}). \end{aligned}$$

Therefore the sample population after one update cycle correctly represents the forward message at time  $t + 1$ .

Particle filtering is *consistent*, therefore, but is it *efficient*? In practice, it seems that the answer is yes: particle filtering seems to maintain a good approximation to the true posterior using a constant number of samples. Under certain assumptions—in particular, that the probabilities in the transition and sensor models are strictly greater than 0 and less than 1—it is possible to prove that the approximation maintains bounded error with high probability. On the practical side, the range of applications has grown to include many fields of science and engineering; some references are given at the end of the chapter.

## 15.6 KEEPING TRACK OF MANY OBJECTS

The preceding sections have considered—without mentioning it—state estimation problems involving a single object. In this section, we see what happens when two or more objects generate the observations. What makes this case different from plain old state estimation is that there is now the possibility of *uncertainty* about which object generated which observation. This is the **identity uncertainty** problem of Section 14.6.3 (page 544), now viewed in a temporal context. In the control theory literature, this is the **data association** problem—that is, the problem of associating observation data with the objects that generated them.



The data association problem was studied originally in the context of radar tracking, where reflected pulses are detected at fixed time intervals by a rotating radar antenna. At each time step, multiple blips may appear on the screen, but there is no direct observation of which blips at time  $t$  belong to which blips at time  $t - 1$ . Figure 15.19(a) shows a simple example with two blips per time step for five steps. Let the two blip locations at time  $t$  be  $e_t^1$  and  $e_t^2$ . (The labeling of blips within a time step as “1” and “2” is completely arbitrary and carries no information.) Let us assume, for the time being, that exactly two aircraft,  $A$  and  $B$ , generated the blips; their true positions are  $X_t^A$  and  $X_t^B$ . Just to keep things simple, we’ll also assume that each aircraft moves independently according to a known transition model—e.g., a linear Gaussian model as used in the Kalman filter (Section 15.4).

Suppose we try to write down the overall probability model for this scenario, just as we did for general temporal processes in Equation (15.3) on page 569. As usual, the joint distribution factors into contributions for each time step as follows:

$$P(x_{0:t}^A, x_{0:t}^B, e_{1:t}^1, e_{1:t}^2) = P(x_0^A)P(x_0^B) \prod_{i=1}^t P(x_i^A | x_{i-1}^A)P(x_i^B | x_{i-1}^B) P(e_i^1, e_i^2 | x_i^A, x_i^B). \quad (15.24)$$

We would like to factor the observation term  $P(e_i^1, e_i^2 | x_i^A, x_i^B)$  into a product of two terms, one for each object, but this would require knowing which observation was generated by which object. Instead, we have to sum over all possible ways of associating the observations

with the objects. Some of those ways are shown in Figure 15.19(b–c); in general, for  $n$  objects and  $T$  time steps, there are  $(n!)^T$  ways of doing it—an awfully large number.

Mathematically speaking, the “way of associating the observations with the objects” is a collection of unobserved random variable that identify the source of each observation. We’ll write  $\omega_t$  to denote the one-to-one mapping from objects to observations at time  $t$ , with  $\omega_t(A)$  and  $\omega_t(B)$  denoting the specific observations (1 or 2) that  $\omega_t$  assigns to  $A$  and  $B$ . (For  $n$  objects,  $\omega_t$  will have  $n!$  possible values; here,  $n! = 2$ .) Because the labels “1” and “2” on the observations are assigned arbitrarily, the prior on  $\omega_t$  is uniform and  $\omega_t$  is independent of the states of the objects,  $x_t^A$  and  $x_t^B$ . So we can condition the observation term  $P(e_i^1, e_i^2 | x_i^A, x_i^B)$  on  $\omega_i$  and then simplify:

$$\begin{aligned} P(e_i^1, e_i^2 | x_i^A, x_i^B) &= \sum_{\omega_i} P(e_i^1, e_i^2 | x_i^A, x_i^B, \omega_i) P(\omega_i | x_i^A, x_i^B) \\ &= \sum_{\omega_i} P(e_i^{\omega_i(A)} | x_i^A) P(e_i^{\omega_i(B)} | x_i^B) P(\omega_i | x_i^A, x_i^B) \\ &= \frac{1}{2} \sum_{\omega_i} P(e_i^{\omega_i(A)} | x_i^A) P(e_i^{\omega_i(B)} | x_i^B). \end{aligned}$$

Plugging this into Equation (15.24), we get an expression that is only in terms of transition and sensor models for individual objects and observations.

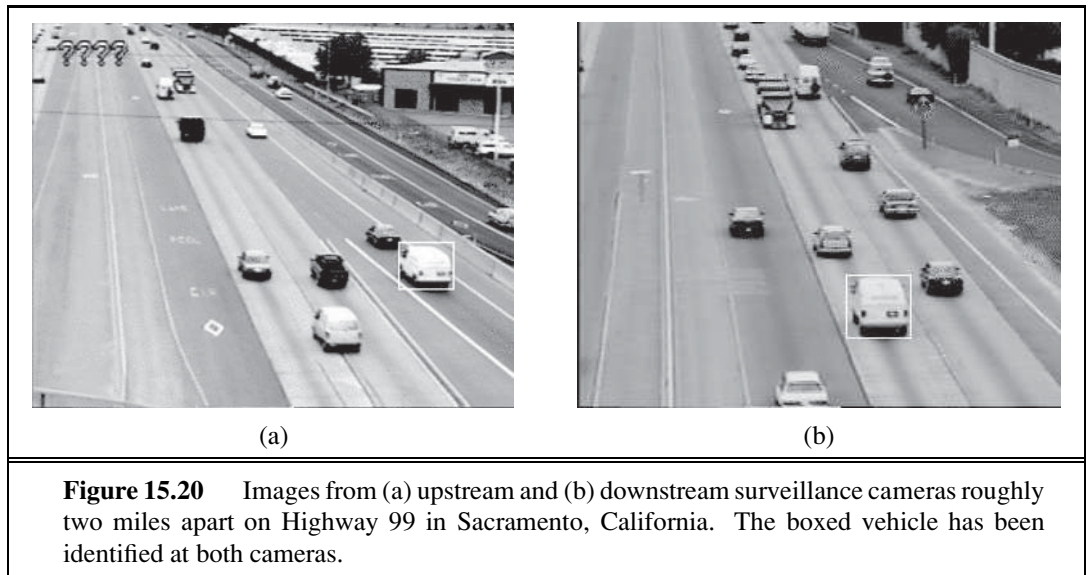
As for all probability models, inference means summing out the variables other than the query and the evidence. For filtering in HMMs and DBNs, we were able to sum out the state variables from 1 to  $t - 1$  by a simple dynamic programming trick; for Kalman filters, we took advantage of special properties of Gaussians. For data association, we are less fortunate. There is no (known) efficient exact algorithm, for the same reason that there is none for the switching Kalman filter (page 589): the filtering distribution  $P(x_t^A | e_{1:t}^1, e_{1:t}^2)$  for object  $A$  ends up as a mixture of exponentially many distributions, one for each way of picking a sequence of observations to assign to  $A$ .

As a result of the complexity of exact inference, many different approximate methods have been used. The simplest approach is to choose a single “best” assignment at each time step, given the predicted positions of the objects at the current time step. This assignment associates observations with objects and enables the track of each object to be updated and a prediction made for the next time step. For choosing the “best” assignment, it is common to use the so-called **nearest-neighbor filter**, which repeatedly chooses the closest pairing of predicted position and observation and adds that pairing to the assignment. The nearest-neighbor filter works well when the objects are well separated in state space and the prediction uncertainty and observation error are small—in other words, when there is no possibility of confusion. When there is more uncertainty as to the correct assignment, a better approach is to choose the assignment that maximizes the joint probability of the current observations given the predicted positions. This can be done very efficiently using the **Hungarian algorithm** (Kuhn, 1955), even though there are  $n!$  assignments to choose from.

Any method that commits to a single best assignment at each time step fails miserably under more difficult conditions. In particular, if the algorithm commits to an incorrect assignment, the prediction at the next time step may be significantly wrong, leading to more

NEAREST-NEIGHBOR  
FILTER

HUNGARIAN  
ALGORITHM



incorrect assignments, and so on. Two modern approaches turn out to be much more effective. A **particle filtering** algorithm (see page 598) for data association works by maintaining a large collection of possible current assignments. An **MCMC** algorithm explores the space of assignment histories—for example, Figure 15.19(b–c) might be states in the MCMC state space—and can change its mind about previous assignment decisions. Current MCMC data association methods can handle many hundreds of objects in real time while giving a good approximation to the true posterior distributions.

The scenario described so far involved  $n$  known objects generating  $n$  observations at each time step. Real application of data association are typically much more complicated. Often, the reported observations include **false alarms** (also known as **clutter**), which are not caused by real objects. **Detection failures** can occur, meaning that no observation is reported for a real object. Finally, new objects arrive and old ones disappear. These phenomena, which create even more possible worlds to worry about, are illustrated in Figure 15.19(d).

Figure 15.20 shows two images from widely separated cameras on a California freeway. In this application, we are interested in two goals: estimating the time it takes, under current traffic conditions, to go from one place to another in the freeway system; and measuring *demand*, i.e., how many vehicles travel between any two points in the system at particular times of the day and on particular days of the week. Both goals require solving the data association problem over a wide area with many cameras and tens of thousands of vehicles per hour. With visual surveillance, false alarms are caused by moving shadows, articulated vehicles, reflections in puddles, etc.; detection failures are caused by occlusion, fog, darkness, and lack of visual contrast; and vehicles are constantly entering and leaving the freeway system. Furthermore, the appearance of any given vehicle can change dramatically between cameras depending on lighting conditions and vehicle pose in the image, and the transition model changes as traffic jams come and go. Despite these problems, modern data association algorithms have been successful in estimating traffic parameters in real-world settings.

FALSE ALARM

CLUTTER

DETECTION FAILURE

Data association is an essential foundation for keeping track of a complex world, because without it there is no way to combine multiple observations of any given object. When objects in the world interact with each other in complex activities, understanding the world requires combining data association with the relational and open-universe probability models of Section 14.6.3. This is currently an active area of research.

---

## 15.7 SUMMARY

---

This chapter has addressed the general problem of representing and reasoning about probabilistic temporal processes. The main points are as follows:

- The changing state of the world is handled by using a set of random variables to represent the state at each point in time.
- Representations can be designed to satisfy the **Markov property**, so that the future is independent of the past given the present. Combined with the assumption that the process is **stationary**—that is, the dynamics do not change over time—this greatly simplifies the representation.
- A temporal probability model can be thought of as containing a **transition model** describing the state evolution and a **sensor model** describing the observation process.
- The principal inference tasks in temporal models are **filtering**, **prediction**, **smoothing**, and computing the **most likely explanation**. Each of these can be achieved using simple, recursive algorithms whose run time is linear in the length of the sequence.
- Three families of temporal models were studied in more depth: **hidden Markov models**, **Kalman filters**, and **dynamic Bayesian networks** (which include the other two as special cases).
- Unless special assumptions are made, as in Kalman filters, exact inference with many state variables is intractable. In practice, the **particle filtering** algorithm seems to be an effective approximation algorithm.
- When trying to keep track of many objects, uncertainty arises as to which observations belong to which objects—the **data association** problem. The number of association hypotheses is typically intractably large, but MCMC and particle filtering algorithms for data association work well in practice.

---

## BIBLIOGRAPHICAL AND HISTORICAL NOTES

Many of the basic ideas for estimating the state of dynamical systems came from the mathematician C. F. Gauss (1809), who formulated a deterministic least-squares algorithm for the problem of estimating orbits from astronomical observations. A. A. Markov (1913) developed what was later called the **Markov assumption** in his analysis of stochastic processes;

he estimated a first-order Markov chain on letters from the text of *Eugene Onegin*. The general theory of Markov chains and their mixing times is covered by Levin *et al.* (2008).

Significant classified work on filtering was done during World War II by Wiener (1942) for continuous-time processes and by Kolmogorov (1941) for discrete-time processes. Although this work led to important technological developments over the next 20 years, its use of a frequency-domain representation made many calculations quite cumbersome. Direct state-space modeling of the stochastic process turned out to be simpler, as shown by Peter Swerling (1959) and Rudolf Kalman (1960). The latter paper described what is now known as the Kalman filter for forward inference in linear systems with Gaussian noise; Kalman's results had, however, been obtained previously by the Danish statistician Thorvold Thiele (1880) and by the Russian mathematician Ruslan Stratonovich (1959), whom Kalman met in Moscow in 1960. After a visit to NASA Ames Research Center in 1960, Kalman saw the applicability of the method to the tracking of rocket trajectories, and the filter was later implemented for the Apollo missions. Important results on smoothing were derived by Rauch *et al.* (1965), and the impressively named Rauch–Tung–Striebel smoother is still a standard technique today. Many early results are gathered in Gelb (1974). Bar-Shalom and Fortmann (1988) give a more modern treatment with a Bayesian flavor, as well as many references to the vast literature on the subject. Chatfield (1989) and Box *et al.* (1994) cover the control theory approach to time series analysis.

The hidden Markov model and associated algorithms for inference and learning, including the forward–backward algorithm, were developed by Baum and Petrie (1966). The Viterbi algorithm first appeared in (Viterbi, 1967). Similar ideas also appeared independently in the Kalman filtering community (Rauch *et al.*, 1965). The forward–backward algorithm was one of the main precursors of the general formulation of the EM algorithm (Dempster *et al.*, 1977); see also Chapter 20. Constant-space smoothing appears in Binder *et al.* (1997b), as does the divide-and-conquer algorithm developed in Exercise 15.3. Constant-time fixed-lag smoothing for HMMs first appeared in Russell and Norvig (2003). HMMs have found many applications in language processing (Charniak, 1993), speech recognition (Rabiner and Juang, 1993), machine translation (Och and Ney, 2003), computational biology (Krogh *et al.*, 1994; Baldi *et al.*, 1994), financial economics Bhar and Hamori (2004) and other fields. There have been several extensions to the basic HMM model, for example the Hierarchical HMM (Fine *et al.*, 1998) and Layered HMM (Oliver *et al.*, 2004) introduce structure back into the model, replacing the single state variable of HMMs.

Dynamic Bayesian networks (DBNs) can be viewed as a sparse encoding of a Markov process and were first used in AI by Dean and Kanazawa (1989b), Nicholson and Brady (1992), and Kjaerulff (1992). The last work extends the HUGIN Bayes net system to accommodate dynamic Bayesian networks. The book by Dean and Wellman (1991) helped popularize DBNs and the probabilistic approach to planning and control within AI. Murphy (2002) provides a thorough analysis of DBNs.

Dynamic Bayesian networks have become popular for modeling a variety of complex motion processes in computer vision (Huang *et al.*, 1994; Intille and Bobick, 1999). Like HMMs, they have found applications in speech recognition (Zweig and Russell, 1998; Richardson *et al.*, 2000; Stephenson *et al.*, 2000; Nefian *et al.*, 2002; Livescu *et al.*, 2003), ge-



nomics (Murphy and Mian, 1999; Perrin *et al.*, 2003; Husmeier, 2003) and robot localization (Theocharous *et al.*, 2004). The link between HMMs and DBNs, and between the forward–backward algorithm and Bayesian network propagation, was made explicitly by Smyth *et al.* (1997). A further unification with Kalman filters (and other statistical models) appears in Roweis and Ghahramani (1999). Procedures exist for learning the parameters (Binder *et al.*, 1997a; Ghahramani, 1998) and structures (Friedman *et al.*, 1998) of DBNs.

The particle filtering algorithm described in Section 15.5 has a particularly interesting history. The first sampling algorithms for particle filtering (also called sequential Monte Carlo methods) were developed in the control theory community by Handschin and Mayne (1969), and the resampling idea that is the core of particle filtering appeared in a Russian control journal (Zaritskii *et al.*, 1975). It was later reinvented in statistics as **sequential importance-sampling resampling**, or **SIR** (Rubin, 1988; Liu and Chen, 1998), in control theory as particle filtering (Gordon *et al.*, 1993; Gordon, 1994), in AI as **survival of the fittest** (Kanazawa *et al.*, 1995), and in computer vision as **condensation** (Isard and Blake, 1996). The paper by Kanazawa *et al.* (1995) includes an improvement called **evidence reversal** whereby the state at time  $t + 1$  is sampled conditional on both the state at time  $t$  and the evidence at time  $t + 1$ . This allows the evidence to influence sample generation directly and was proved by Doucet (1997) and Liu and Chen (1998) to reduce the approximation error. Particle filtering has been applied in many areas, including tracking complex motion patterns in video (Isard and Blake, 1996), predicting the stock market (de Freitas *et al.*, 2000), and diagnosing faults on planetary rovers (Verma *et al.*, 2004). A variant called the **Rao-Blackwellized particle filter** or **RBPF** (Doucet *et al.*, 2000; Murphy and Russell, 2001) applies particle filtering to a subset of state variables and, for each particle, performs exact inference on the remaining variables conditioned on the value sequence in the particle. In some cases RBPF works well with thousands of state variables. An application of RBPF to localization and mapping in robotics is described in Chapter 25. The book by Doucet *et al.* (2001) collects many important papers on **sequential Monte Carlo** (SMC) algorithms, of which particle filtering is the most important instance. Pierre Del Moral and colleagues have performed extensive theoretical analyses of SMC algorithms (Del Moral, 2004; Del Moral *et al.*, 2006).

MCMC methods (see Section 14.5.2) can be applied to the filtering problem; for example, Gibbs sampling can be applied directly to an unrolled DBN. To avoid the problem of increasing update times as the unrolled network grows, the **decayed MCMC** filter (Marthi *et al.*, 2002) prefers to sample more recent state variables, with a probability that decays as  $1/k^2$  for a variable  $k$  steps into the past. Decayed MCMC is a provably nondivergent filter. Nondivergence theorems can also be obtained for certain types of **assumed-density filter**. An assumed-density filter assumes that the posterior distribution over states at time  $t$  belongs to a particular finitely parameterized family; if the projection and update steps take it outside this family, the distribution is projected back to give the best approximation within the family. For DBNs, the Boyen–Koller algorithm (Boyen *et al.*, 1999) and the **factored frontier** algorithm (Murphy and Weiss, 2001) assume that the posterior distribution can be approximated well by a product of small factors. Variational techniques (see Chapter 14) have also been developed for temporal models. Ghahramani and Jordan (1997) discuss an approximation algorithm for the **factorial HMM**, a DBN in which two or more independently evolving

EVIDENCE  
REVERSALRAO-  
BLACKWELLIZED  
PARTICLE FILTERSEQUENTIAL MONTE  
CARLO

DECAYED MCMC

ASSUMED-DENSITY  
FILTERFACTORED  
FRONTIER

FACTORIAL HMM

Markov chains are linked by a shared observation stream. Jordan *et al.* (1998) cover a number of other applications.

Data association for multitarget tracking was first described in a probabilistic setting by Sittler (1964). The first practical algorithm for large-scale problems was the “multiple hypothesis tracker” or MHT algorithm (Reid, 1979). Many important papers are collected by Bar-Shalom and Fortmann (1988) and Bar-Shalom (1992). The development of an MCMC algorithm for data association is due to Pasula *et al.* (1999), who applied it to traffic surveillance problems. Oh *et al.* (2009) provide a formal analysis and extensive experimental comparisons to other methods. Schulz *et al.* (2003) describe a data association method based on particle filtering. Ingemar Cox analyzed the complexity of data association (Cox, 1993; Cox and Hingorani, 1994) and brought the topic to the attention of the vision community. He also noted the applicability of the polynomial-time Hungarian algorithm to the problem of finding most-likely assignments, which had long been considered an intractable problem in the tracking community. The algorithm itself was published by Kuhn (1955), based on translations of papers published in 1931 by two Hungarian mathematicians, Dénes König and Jenő Egerváry. The basic theorem had been derived previously, however, in an unpublished Latin manuscript by the famous Prussian mathematician Carl Gustav Jacobi (1804–1851).

---

## EXERCISES

**15.1** Show that any second-order Markov process can be rewritten as a first-order Markov process with an augmented set of state variables. Can this always be done *parsimoniously*, i.e., without increasing the number of parameters needed to specify the transition model?

**15.2** In this exercise, we examine what happens to the probabilities in the umbrella world in the limit of long time sequences.

- a. Suppose we observe an unending sequence of days on which the umbrella appears. Show that, as the days go by, the probability of rain on the current day increases monotonically toward a fixed point. Calculate this fixed point.
- b. Now consider *forecasting* further and further into the future, given just the first two umbrella observations. First, compute the probability  $P(r_{2+k}|u_1, u_2)$  for  $k = 1 \dots 20$  and plot the results. You should see that the probability converges towards a fixed point. Prove that the exact value of this fixed point is 0.5.

**15.3** This exercise develops a space-efficient variant of the forward-backward algorithm described in Figure 15.4 (page 576). We wish to compute  $\mathbf{P}(\mathbf{X}_k|\mathbf{e}_{1:t})$  for  $k = 1, \dots, t$ . This will be done with a divide-and-conquer approach.

- a. Suppose, for simplicity, that  $t$  is odd, and let the halfway point be  $h = (t + 1)/2$ . Show that  $\mathbf{P}(\mathbf{X}_k|\mathbf{e}_{1:t})$  can be computed for  $k = 1, \dots, h$  given just the initial forward message  $\mathbf{f}_{1:0}$ , the backward message  $\mathbf{b}_{h+1:t}$ , and the evidence  $\mathbf{e}_{1:h}$ .
- b. Show a similar result for the second half of the sequence.

- c. Given the results of (a) and (b), a recursive divide-and-conquer algorithm can be constructed by first running forward along the sequence and then backward from the end, storing just the required messages at the middle and the ends. Then the algorithm is called on each half. Write out the algorithm in detail.
- d. Compute the time and space complexity of the algorithm as a function of  $t$ , the length of the sequence. How does this change if we divide the input into more than two pieces?

**15.4** On page 577, we outlined a flawed procedure for finding the most likely state sequence, given an observation sequence. The procedure involves finding the most likely state at each time step, using smoothing, and returning the sequence composed of these states. Show that, for some temporal probability models and observation sequences, this procedure returns an impossible state sequence (i.e., the posterior probability of the sequence is zero).

**15.5** Equation (15.12) describes the filtering process for the matrix formulation of HMMs. Give a similar equation for the calculation of likelihoods, which was described generically in Equation (15.7).

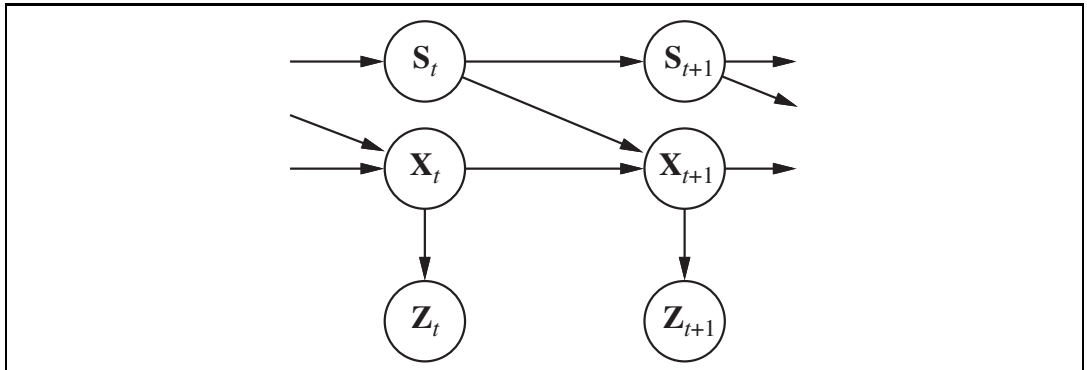
**15.6** Consider the vacuum worlds of Figure 4.18 (perfect sensing) and Figure 15.7 (noisy sensing). Suppose that the robot receives an observation sequence such that, with perfect sensing, there is exactly one possible location it could be in. Is this location necessarily the most probable location under noisy sensing for sufficiently small noise probability  $\epsilon$ ? Prove your claim or find a counterexample.



**15.7** In Section 15.3.2, the prior distribution over locations is uniform and the transition model assumes an equal probability of moving to any neighboring square. What if those assumptions are wrong? Suppose that the initial location is actually chosen uniformly from the northwest quadrant of the room and the *Move* action actually tends to move southeast. Keeping the HMM model fixed, explore the effect on localization and path accuracy as the southeasterly tendency increases, for different values of  $\epsilon$ .

**15.8** Consider a version of the vacuum robot (page 582) that has the policy of going straight for as long as it can; only when it encounters an obstacle does it change to a new (randomly selected) heading. To model this robot, each state in the model consists of a (*location, heading*) pair. Implement this model and see how well the Viterbi algorithm can track a robot with this model. The robot's policy is more constrained than the random-walk robot; does that mean that predictions of the most likely path are more accurate?

**15.9** This exercise is concerned with filtering in an environment with no landmarks. Consider a vacuum robot in an empty room, represented by an  $n \times m$  rectangular grid. The robot's location is hidden; the only evidence available to the observer is a noisy location sensor that gives an approximation to the robot's location. If the robot is at location  $(x, y)$  then with probability .1 the sensor gives the correct location, with probability .05 each it reports one of the 8 locations immediately surrounding  $(x, y)$ , with probability .025 each it reports one of the 16 locations that surround those 8, and with the remaining probability of .1 it reports "no reading." The robot's policy is to pick a direction and follow it with probability .8 on each step; the robot switches to a randomly selected new heading with probability .2 (or with



**Figure 15.21** A Bayesian network representation of a switching Kalman filter. The switching variable  $S_t$  is a discrete state variable whose value determines the transition model for the continuous state variables  $\mathbf{X}_t$ . For any discrete state  $i$ , the transition model  $\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{X}_t, S_t = i)$  is a linear Gaussian model, just as in a regular Kalman filter. The transition model for the discrete state,  $\mathbf{P}(S_{t+1}|S_t)$ , can be thought of as a matrix, as in a hidden Markov model.

probability 1 if it encounters a wall). Implement this as an HMM and do filtering to track the robot. How accurately can we track the robot's path?

**15.10** Often, we wish to monitor a continuous-state system whose behavior switches unpredictably among a set of  $k$  distinct “modes.” For example, an aircraft trying to evade a missile can execute a series of distinct maneuvers that the missile may attempt to track. A Bayesian network representation of such a **switching Kalman filter** model is shown in Figure 15.21.

- Suppose that the discrete state  $S_t$  has  $k$  possible values and that the prior continuous state estimate  $\mathbf{P}(\mathbf{X}_0)$  is a multivariate Gaussian distribution. Show that the prediction  $\mathbf{P}(\mathbf{X}_1)$  is a **mixture of Gaussians**—that is, a weighted sum of Gaussians such that the weights sum to 1.
- Show that if the current continuous state estimate  $\mathbf{P}(\mathbf{X}_t|\mathbf{e}_{1:t})$  is a mixture of  $m$  Gaussians, then in the general case the updated state estimate  $\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1})$  will be a mixture of  $km$  Gaussians.
- What aspect of the temporal process do the weights in the Gaussian mixture represent?

The results in (a) and (b) show that the representation of the posterior grows without limit even for switching Kalman filters, which are among the simplest hybrid dynamic models.

**15.11** Complete the missing step in the derivation of Equation (15.19) on page 586, the first update step for the one-dimensional Kalman filter.

**15.12** Let us examine the behavior of the variance update in Equation (15.20) (page 587).

- Plot the value of  $\sigma_t^2$  as a function of  $t$ , given various values for  $\sigma_x^2$  and  $\sigma_z^2$ .
- Show that the update has a fixed point  $\sigma^2$  such that  $\sigma_t^2 \rightarrow \sigma^2$  as  $t \rightarrow \infty$ , and calculate the value of  $\sigma^2$ .
- Give a qualitative explanation for what happens as  $\sigma_x^2 \rightarrow 0$  and as  $\sigma_z^2 \rightarrow 0$ .

**15.13** A professor wants to know if students are getting enough sleep. Each day, the professor observes whether the students sleep in class, and whether they have red eyes. The professor has the following domain theory:

- The prior probability of getting enough sleep, with no observations, is 0.7.
- The probability of getting enough sleep on night  $t$  is 0.8 given that the student got enough sleep the previous night, and 0.3 if not.
- The probability of having red eyes is 0.2 if the student got enough sleep, and 0.7 if not.
- The probability of sleeping in class is 0.1 if the student got enough sleep, and 0.3 if not.

Formulate this information as a dynamic Bayesian network that the professor could use to filter or predict from a sequence of observations. Then reformulate it as a hidden Markov model that has only a single observation variable. Give the complete probability tables for the model.

**15.14** For the DBN specified in Exercise 15.13 and for the evidence values

- $\mathbf{e}_1$  = not red eyes, not sleeping in class
- $\mathbf{e}_2$  = red eyes, not sleeping in class
- $\mathbf{e}_3$  = red eyes, sleeping in class

perform the following computations:

- a. State estimation: Compute  $P(\text{EnoughSleep}_t | \mathbf{e}_{1:t})$  for each of  $t = 1, 2, 3$ .
- b. Smoothing: Compute  $P(\text{EnoughSleep}_t | \mathbf{e}_{1:3})$  for each of  $t = 1, 2, 3$ .
- c. Compare the filtered and smoothed probabilities for  $t = 1$  and  $t = 2$ .

**15.15** Suppose that a particular student shows up with red eyes and sleeps in class every day. Given the model described in Exercise 15.13, explain why the probability that the student had enough sleep the previous night converges to a fixed point rather than continuing to go down as we gather more days of evidence. What is the fixed point? Answer this both numerically (by computation) and analytically.

**15.16** This exercise analyzes in more detail the persistent-failure model for the battery sensor in Figure 15.15(a) (page 594).

- a. Figure 15.15(b) stops at  $t = 32$ . Describe qualitatively what should happen as  $t \rightarrow \infty$  if the sensor continues to read 0.
- b. Suppose that the external temperature affects the battery sensor in such a way that transient failures become more likely as temperature increases. Show how to augment the DBN structure in Figure 15.15(a), and explain any required changes to the CPTs.
- c. Given the new network structure, can battery readings be used by the robot to infer the current temperature?

**15.17** Consider applying the variable elimination algorithm to the umbrella DBN unrolled for three slices, where the query is  $\mathbf{P}(R_3 | u_1, u_2, u_3)$ . Show that the space complexity of the algorithm—the size of the largest factor—is the same, regardless of whether the rain variables are eliminated in forward or backward order.

# 16 MAKING SIMPLE DECISIONS

*In which we see how an agent should make decisions so that it gets what it wants—on average, at least.*

In this chapter, we fill in the details of how utility theory combines with probability theory to yield a decision-theoretic agent—an agent that can make rational decisions based on what it believes and what it wants. Such an agent can make decisions in contexts in which uncertainty and conflicting goals leave a logical agent with no way to decide: a goal-based agent has a binary distinction between good (goal) and bad (non-goal) states, while a decision-theoretic agent has a continuous measure of outcome quality.

Section 16.1 introduces the basic principle of decision theory: the maximization of expected utility. Section 16.2 shows that the behavior of any rational agent can be captured by supposing a utility function that is being maximized. Section 16.3 discusses the nature of utility functions in more detail, and in particular their relation to individual quantities such as money. Section 16.4 shows how to handle utility functions that depend on several quantities. In Section 16.5, we describe the implementation of decision-making systems. In particular, we introduce a formalism called a **decision network** (also known as an **influence diagram**) that extends Bayesian networks by incorporating actions and utilities. The remainder of the chapter discusses issues that arise in applications of decision theory to expert systems.

## 16.1 COMBINING BELIEFS AND DESIRES UNDER UNCERTAINTY

Decision theory, in its simplest form, deals with choosing among actions based on the desirability of their *immediate* outcomes; that is, the environment is assumed to be episodic in the sense defined on page 43. (This assumption is relaxed in Chapter 17.) In Chapter 3 we used the notation  $\text{RESULT}(s_0, a)$  for the state that is the deterministic outcome of taking action  $a$  in state  $s_0$ . In this chapter we deal with nondeterministic partially observable environments. Since the agent may not know the current state, we omit it and define  $\text{RESULT}(a)$  as a *random variable* whose values are the possible outcome states. The probability of outcome  $s'$ , given evidence observations  $\mathbf{e}$ , is written

$$P(\text{RESULT}(a) = s' \mid a, \mathbf{e}),$$

where the  $a$  on the right-hand side of the conditioning bar stands for the event that action  $a$  is executed.<sup>1</sup>

UTILITY FUNCTION

EXPECTED UTILITY

The agent's preferences are captured by a **utility function**,  $U(s)$ , which assigns a single number to express the desirability of a state. The **expected utility** of an action given the evidence,  $EU(a|\mathbf{e})$ , is just the average utility value of the outcomes, weighted by the probability that the outcome occurs:

$$EU(a|\mathbf{e}) = \sum_{s'} P(\text{RESULT}(a) = s' | a, \mathbf{e}) U(s'). \quad (16.1)$$

MAXIMUM EXPECTED UTILITY

The principle of **maximum expected utility** (MEU) says that a rational agent should choose the action that maximizes the agent's expected utility:

$$\text{action} = \underset{a}{\operatorname{argmax}} EU(a|\mathbf{e})$$

In a sense, the MEU principle could be seen as defining all of AI. All an intelligent agent has to do is calculate the various quantities, maximize utility over its actions, and away it goes. But this does not mean that the AI problem is *solved* by the definition!

The MEU principle *formalizes* the general notion that the agent should “do the right thing,” but goes only a small distance toward a full *operationalization* of that advice. Estimating the state of the world requires perception, learning, knowledge representation, and inference. Computing  $P(\text{RESULT}(a) | a, \mathbf{e})$  requires a complete causal model of the world and, as we saw in Chapter 14, NP-hard inference in (very large) Bayesian networks. Computing the outcome utilities  $U(s')$  often requires searching or planning, because an agent may not know how good a state is until it knows where it can get to from that state. So, decision theory is not a panacea that solves the AI problem—but it does provide a useful framework.

The MEU principle has a clear relation to the idea of performance measures introduced in Chapter 2. The basic idea is simple. Consider the environments that could lead to an agent having a given percept history, and consider the different agents that we could design. *If an agent acts so as to maximize a utility function that correctly reflects the performance measure, then the agent will achieve the highest possible performance score (averaged over all the possible environments).* This is the central justification for the MEU principle itself. While the claim may seem tautological, it does in fact embody a very important transition from a global, external criterion of rationality—the performance measure over environment histories—to a local, internal criterion involving the maximization of a utility function applied to the next state.



## 16.2 THE BASIS OF UTILITY THEORY

Intuitively, the principle of Maximum Expected Utility (MEU) seems like a reasonable way to make decisions, but it is by no means obvious that it is the *only* rational way. After all, why should maximizing the *average* utility be so special? What's wrong with an agent that

<sup>1</sup> Classical decision theory leaves the current state  $S_0$  implicit, but we could make it explicit by writing  $P(\text{RESULT}(a) = s' | a, \mathbf{e}) = \sum_s P(\text{RESULT}(s, a) = s' | a) P(S_0 = s | \mathbf{e})$ .

maximizes the weighted sum of the cubes of the possible utilities, or tries to minimize the worst possible loss? Could an agent act rationally just by expressing preferences between states, without giving them numeric values? Finally, why should a utility function with the required properties exist at all? We shall see.

### 16.2.1 Constraints on rational preferences

These questions can be answered by writing down some constraints on the preferences that a rational agent should have and then showing that the MEU principle can be derived from the constraints. We use the following notation to describe an agent's preferences:

$A \succ B$  the agent prefers  $A$  over  $B$ .

$A \sim B$  the agent is indifferent between  $A$  and  $B$ .

$A \succsim B$  the agent prefers  $A$  over  $B$  or is indifferent between them.

Now the obvious question is, what sorts of things are  $A$  and  $B$ ? They could be states of the world, but more often than not there is uncertainty about what is really being offered. For example, an airline passenger who is offered “the pasta dish or the chicken” does not know what lurks beneath the tinfoil cover.<sup>2</sup> The pasta could be delicious or congealed, the chicken juicy or overcooked beyond recognition. We can think of the set of outcomes for each action as a **lottery**—think of each action as a ticket. A lottery  $L$  with possible outcomes  $S_1, \dots, S_n$  that occur with probabilities  $p_1, \dots, p_n$  is written

$$L = [p_1, S_1; p_2, S_2; \dots p_n, S_n].$$

In general, each outcome  $S_i$  of a lottery can be either an atomic state or another lottery. The primary issue for utility theory is to understand how preferences between complex lotteries are related to preferences between the underlying states in those lotteries. To address this issue we list six constraints that we require any reasonable preference relation to obey:

LOTTERY

ORDERABILITY

- **Orderability:** Given any two lotteries, a rational agent must either prefer one to the other or else rate the two as equally preferable. That is, the agent cannot avoid deciding. As we said on page 490, refusing to bet is like refusing to allow time to pass.

Exactly one of  $(A \succ B)$ ,  $(B \succ A)$ , or  $(A \sim B)$  holds.

TRANSITIVITY

- **Transitivity:** Given any three lotteries, if an agent prefers  $A$  to  $B$  and prefers  $B$  to  $C$ , then the agent must prefer  $A$  to  $C$ .

$$(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C).$$

CONTINUITY

- **Continuity:** If some lottery  $B$  is between  $A$  and  $C$  in preference, then there is some probability  $p$  for which the rational agent will be indifferent between getting  $B$  for sure and the lottery that yields  $A$  with probability  $p$  and  $C$  with probability  $1 - p$ .

$$A \succ B \succ C \Rightarrow \exists p [p, A; 1 - p, C] \sim B.$$

SUBSTITUTABILITY

- **Substitutability:** If an agent is indifferent between two lotteries  $A$  and  $B$ , then the agent is indifferent between two more complex lotteries that are the same except that  $B$

<sup>2</sup> We apologize to readers whose local airlines no longer offer food on long flights.



is substituted for  $A$  in one of them. This holds regardless of the probabilities and the other outcome(s) in the lotteries.

$$A \sim B \Rightarrow [p, A; 1 - p, C] \sim [p, B; 1 - p, C] .$$

This also holds if we substitute  $\succ$  for  $\sim$  in this axiom.

MONOTONICITY

- **Monotonicity:** Suppose two lotteries have the same two possible outcomes,  $A$  and  $B$ . If an agent prefers  $A$  to  $B$ , then the agent must prefer the lottery that has a higher probability for  $A$  (and vice versa).

$$A \succ B \Rightarrow (p > q \Leftrightarrow [p, A; 1 - p, B] \succ [q, A; 1 - q, B]) .$$

DECOMPOSABILITY

- **Decomposability:** Compound lotteries can be reduced to simpler ones using the laws of probability. This has been called the “no fun in gambling” rule because it says that two consecutive lotteries can be compressed into a single equivalent lottery, as shown in Figure 16.1(b).<sup>3</sup>

$$[p, A; 1 - p, [q, B; 1 - q, C]] \sim [p, A; (1 - p)q, B; (1 - p)(1 - q), C] .$$

These constraints are known as the axioms of utility theory. Each axiom can be motivated by showing that an agent that violates it will exhibit patently irrational behavior in some situations. For example, we can motivate transitivity by making an agent with nontransitive preferences give us all its money. Suppose that the agent has the nontransitive preferences  $A \succ B \succ C \succ A$ , where  $A$ ,  $B$ , and  $C$  are goods that can be freely exchanged. If the agent currently has  $A$ , then we could offer to trade  $C$  for  $A$  plus one cent. The agent prefers  $C$ , and so would be willing to make this trade. We could then offer to trade  $B$  for  $C$ , extracting another cent, and finally trade  $A$  for  $B$ . This brings us back where we started from, except that the agent has given us three cents (Figure 16.1(a)). We can keep going around the cycle until the agent has no money at all. Clearly, the agent has acted irrationally in this case.

### 16.2.2 Preferences lead to utility

Notice that the axioms of utility theory are really axioms about preferences—they say nothing about a utility function. But in fact from the axioms of utility we can derive the following consequences (for the proof, see von Neumann and Morgenstern, 1944):

- **Existence of Utility Function:** If an agent’s preferences obey the axioms of utility, then there exists a function  $U$  such that  $U(A) > U(B)$  if and only if  $A$  is preferred to  $B$ , and  $U(A) = U(B)$  if and only if the agent is indifferent between  $A$  and  $B$ .

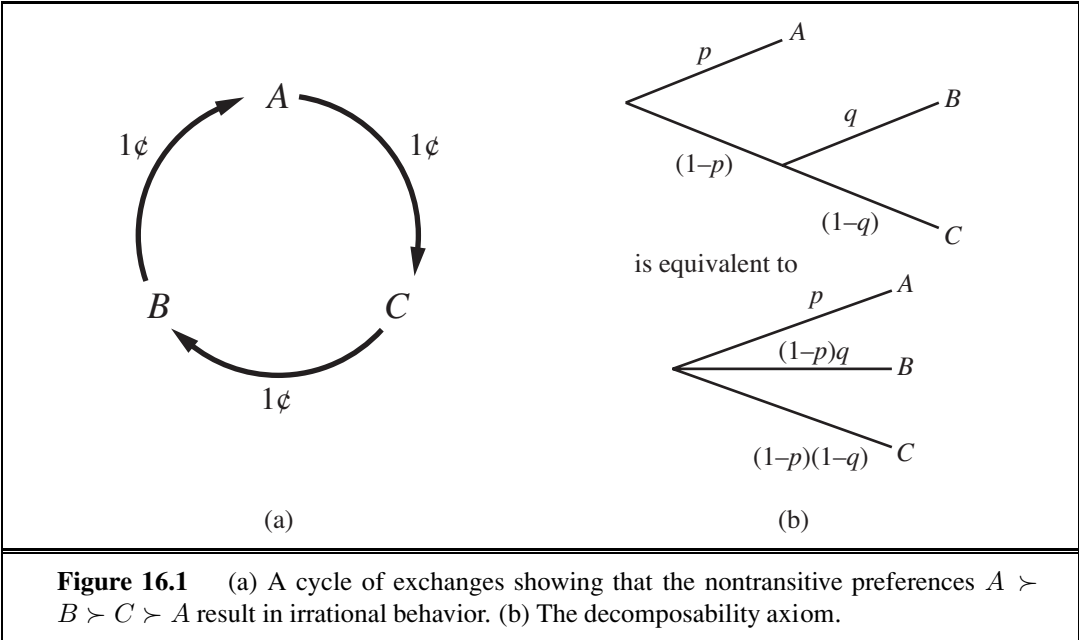
$$U(A) > U(B) \Leftrightarrow A \succ B$$

$$U(A) = U(B) \Leftrightarrow A \sim B$$

- **Expected Utility of a Lottery:** The utility of a lottery is the sum of the probability of each outcome times the utility of that outcome.

$$U([p_1, S_1; \dots; p_n, S_n]) = \sum_i p_i U(S_i) .$$

<sup>3</sup> We can account for the enjoyment of gambling by encoding gambling events into the state description; for example, “Have \$10 and gambled” could be preferred to “Have \$10 and didn’t gamble.”



In other words, once the probabilities and utilities of the possible outcome states are specified, the utility of a compound lottery involving those states is completely determined. Because the outcome of a nondeterministic action is a lottery, it follows that an agent can act rationally—that is, consistently with its preferences—only by choosing an action that maximizes expected utility according to Equation (16.1).

The preceding theorems establish that a utility function *exists* for any rational agent, but they do not establish that it is *unique*. It is easy to see, in fact, that an agent’s behavior would not change if its utility function  $U(S)$  were transformed according to

$$U'(S) = aU(S) + b,$$

(16.2)

where  $a$  and  $b$  are constants and  $a > 0$ ; an affine transformation.<sup>4</sup> This fact was noted in Chapter 5 for two-player games of chance; here, we see that it is completely general.

As in game-playing, in a deterministic environment an agent just needs a preference ranking on states—the numbers don’t matter. This is called a **value function** or **ordinal utility function**.

It is important to remember that the existence of a utility function that describes an agent’s preference behavior does not necessarily mean that the agent is *explicitly* maximizing that utility function in its own deliberations. As we showed in Chapter 2, rational behavior can be generated in any number of ways. By observing a rational agent’s preferences, however, an observer can construct the utility function that represents what the agent is actually trying to achieve (even if the agent doesn’t know it).

<sup>4</sup> In this sense, utilities resemble temperatures: a temperature in Fahrenheit is 1.8 times the Celsius temperature plus 32. You get the same results in either measurement system.

VALUE FUNCTION  
ORDINAL UTILITY  
FUNCTION

## 16.3 UTILITY FUNCTIONS

Utility is a function that maps from lotteries to real numbers. We know there are some axioms on utilities that all rational agents must obey. Is that all we can say about utility functions? Strictly speaking, that is it: an agent can have any preferences it likes. For example, an agent might prefer to have a prime number of dollars in its bank account; in which case, if it had \$16 it would give away \$3. This might be unusual, but we can't call it irrational. An agent might prefer a dented 1973 Ford Pinto to a shiny new Mercedes. Preferences can also interact: for example, the agent might prefer prime numbers of dollars only when it owns the Pinto, but when it owns the Mercedes, it might prefer more dollars to fewer. Fortunately, the preferences of real agents are usually more systematic, and thus easier to deal with.

### 16.3.1 Utility assessment and utility scales

If we want to build a decision-theoretic system that helps the agent make decisions or acts on his or her behalf, we must first work out what the agent's utility function is. This process, often called **preference elicitation**, involves presenting choices to the agent and using the observed preferences to pin down the underlying utility function.

PREFERENCE  
ELICITATION

Equation (16.2) says that there is no absolute scale for utilities, but it is helpful, nonetheless, to establish *some* scale on which utilities can be recorded and compared for any particular problem. A scale can be established by fixing the utilities of any two particular outcomes, just as we fix a temperature scale by fixing the freezing point and boiling point of water. Typically, we fix the utility of a “best possible prize” at  $U(S) = u_{\top}$  and a “worst possible catastrophe” at  $U(S) = u_{\perp}$ . **Normalized utilities** use a scale with  $u_{\perp} = 0$  and  $u_{\top} = 1$ .

NORMALIZED  
UTILITIES

Given a utility scale between  $u_{\top}$  and  $u_{\perp}$ , we can assess the utility of any particular prize  $S$  by asking the agent to choose between  $S$  and a **standard lottery**  $[p, u_{\top}; (1-p), u_{\perp}]$ . The probability  $p$  is adjusted until the agent is indifferent between  $S$  and the standard lottery. Assuming normalized utilities, the utility of  $S$  is given by  $p$ . Once this is done for each prize, the utilities for all lotteries involving those prizes are determined.

STANDARD LOTTERY



In medical, transportation, and environmental decision problems, among others, people's lives are at stake. In such cases,  $u_{\perp}$  is the value assigned to immediate death (or perhaps many deaths). *Although nobody feels comfortable with putting a value on human life, it is a fact that tradeoffs are made all the time.* Aircraft are given a complete overhaul at intervals determined by trips and miles flown, rather than after every trip. Cars are manufactured in a way that trades off costs against accident survival rates. Paradoxically, a refusal to “put a monetary value on life” means that life is often *undervalued*. Ross Shachter relates an experience with a government agency that commissioned a study on removing asbestos from schools. The decision analysts performing the study assumed a particular dollar value for the life of a school-age child, and argued that the rational choice under that assumption was to remove the asbestos. The agency, morally outraged at the idea of setting the value of a life, rejected the report out of hand. It then decided against asbestos removal—implicitly asserting a lower value for the life of a child than that assigned by the analysts.

MICROMORT

Some attempts have been made to find out the value that people place on their own lives. One common “currency” used in medical and safety analysis is the **micromort**, a one in a million chance of death. If you ask people how much they would pay to avoid a risk—for example, to avoid playing Russian roulette with a million-barreled revolver—they will respond with very large numbers, perhaps tens of thousands of dollars, but their actual behavior reflects a much lower monetary value for a micromort. For example, driving in a car for 230 miles incurs a risk of one micromort; over the life of your car—say, 92,000 miles—that’s 400 micromorts. People appear to be willing to pay about \$10,000 (at 2009 prices) more for a safer car that halves the risk of death, or about \$50 per micromort. A number of studies have confirmed a figure in this range across many individuals and risk types. Of course, this argument holds only for small risks. Most people won’t agree to kill themselves for \$50 million.

QALY

Another measure is the **QALY**, or quality-adjusted life year. Patients with a disability are willing to accept a shorter life expectancy to be restored to full health. For example, kidney patients on average are indifferent between living two years on a dialysis machine and one year at full health.

### 16.3.2 The utility of money

Utility theory has its roots in economics, and economics provides one obvious candidate for a utility measure: money (or more specifically, an agent’s total net assets). The almost universal exchangeability of money for all kinds of goods and services suggests that money plays a significant role in human utility functions.

MONOTONIC  
PREFERENCE

It will usually be the case that an agent prefers more money to less, all other things being equal. We say that the agent exhibits a **monotonic preference** for more money. This does not mean that money behaves as a utility function, because it says nothing about preferences between *lotteries* involving money.

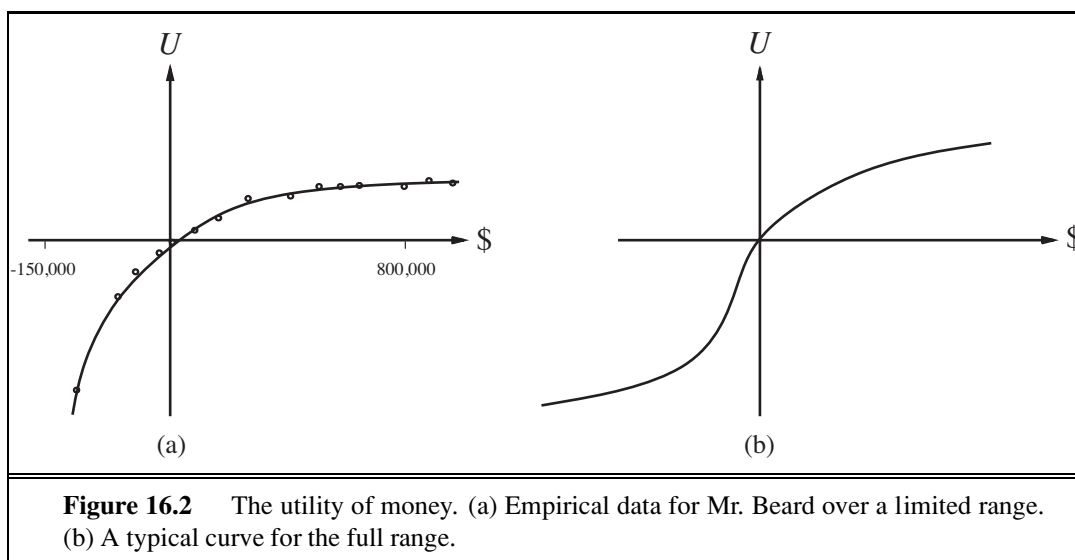
Suppose you have triumphed over the other competitors in a television game show. The host now offers you a choice: either you can take the \$1,000,000 prize or you can gamble it on the flip of a coin. If the coin comes up heads, you end up with nothing, but if it comes up tails, you get \$2,500,000. If you’re like most people, you would decline the gamble and pocket the million. Are you being irrational?

EXPECTED  
MONETARY VALUE

Assuming the coin is fair, the **expected monetary value** (EMV) of the gamble is  $\frac{1}{2}(\$0) + \frac{1}{2}(\$2,500,000) = \$1,250,000$ , which is more than the original \$1,000,000. But that does not necessarily mean that accepting the gamble is a better decision. Suppose we use  $S_n$  to denote the state of possessing total wealth  $\$n$ , and that your current wealth is  $\$k$ . Then the expected utilities of the two actions of accepting and declining the gamble are

$$\begin{aligned} EU(\text{Accept}) &= \frac{1}{2}U(S_k) + \frac{1}{2}U(S_{k+2,500,000}) , \\ EU(\text{Decline}) &= U(S_{k+1,000,000}) . \end{aligned}$$

To determine what to do, we need to assign utilities to the outcome states. Utility is not directly proportional to monetary value, because the utility for your first million is very high (or so they say), whereas the utility for an additional million is smaller. Suppose you assign a utility of 5 to your current financial status ( $S_k$ ), a 9 to the state  $S_{k+2,500,000}$ , and an 8 to the



state  $S_{k+1,000,000}$ . Then the rational action would be to decline, because the expected utility of accepting is only 7 (less than the 8 for declining). On the other hand, a billionaire would most likely have a utility function that is locally linear over the range of a few million more, and thus would accept the gamble.

In a pioneering study of actual utility functions, Grayson (1960) found that the utility of money was almost exactly proportional to the *logarithm* of the amount. (This idea was first suggested by Bernoulli (1738); see Exercise 16.3.) One particular utility curve, for a certain Mr. Beard, is shown in Figure 16.2(a). The data obtained for Mr. Beard's preferences are consistent with a utility function

$$U(S_{k+n}) = -263.31 + 22.09 \log(n + 150,000)$$

for the range between  $n = -\$150,000$  and  $n = \$800,000$ .

We should not assume that this is the definitive utility function for monetary value, but it is likely that most people have a utility function that is concave for positive wealth. Going into debt is bad, but preferences between different levels of debt can display a reversal of the concavity associated with positive wealth. For example, someone already \$10,000,000 in debt might well accept a gamble on a fair coin with a gain of \$10,000,000 for heads and a loss of \$20,000,000 for tails.<sup>5</sup> This yields the S-shaped curve shown in Figure 16.2(b).

If we restrict our attention to the positive part of the curves, where the slope is decreasing, then for any lottery  $L$ , the utility of being faced with that lottery is less than the utility of being handed the expected monetary value of the lottery as a sure thing:

$$U(L) < U(S_{EMV(L)}).$$

RISK-AVERSE

That is, agents with curves of this shape are **risk-averse**: they prefer a sure thing with a payoff that is less than the expected monetary value of a gamble. On the other hand, in the “desperate” region at large negative wealth in Figure 16.2(b), the behavior is **risk-seeking**.

RISK-SEEKING

<sup>5</sup> Such behavior might be called desperate, but it is rational if one is already in a desperate situation.

CERTAINTY  
EQUIVALENT

The value an agent will accept in lieu of a lottery is called the **certainty equivalent** of the lottery. Studies have shown that most people will accept about \$400 in lieu of a gamble that gives \$1000 half the time and \$0 the other half—that is, the certainty equivalent of the lottery is \$400, while the EMV is \$500. The difference between the EMV of a lottery and its certainty equivalent is called the **insurance premium**. Risk aversion is the basis for the insurance industry, because it means that insurance premiums are positive. People would rather pay a small insurance premium than gamble the price of their house against the chance of a fire. From the insurance company’s point of view, the price of the house is very small compared with the firm’s total reserves. This means that the insurer’s utility curve is approximately linear over such a small region, and the gamble costs the company almost nothing.

INSURANCE  
PREMIUM

RISK-NEUTRAL

Notice that for *small* changes in wealth relative to the current wealth, almost any curve will be approximately linear. An agent that has a linear curve is said to be **risk-neutral**. For gambles with small sums, therefore, we expect risk neutrality. In a sense, this justifies the simplified procedure that proposed small gambles to assess probabilities and to justify the axioms of probability in Section 13.2.3.

### 16.3.3 Expected utility and post-decision disappointment

The rational way to choose the best action,  $a^*$ , is to maximize expected utility:

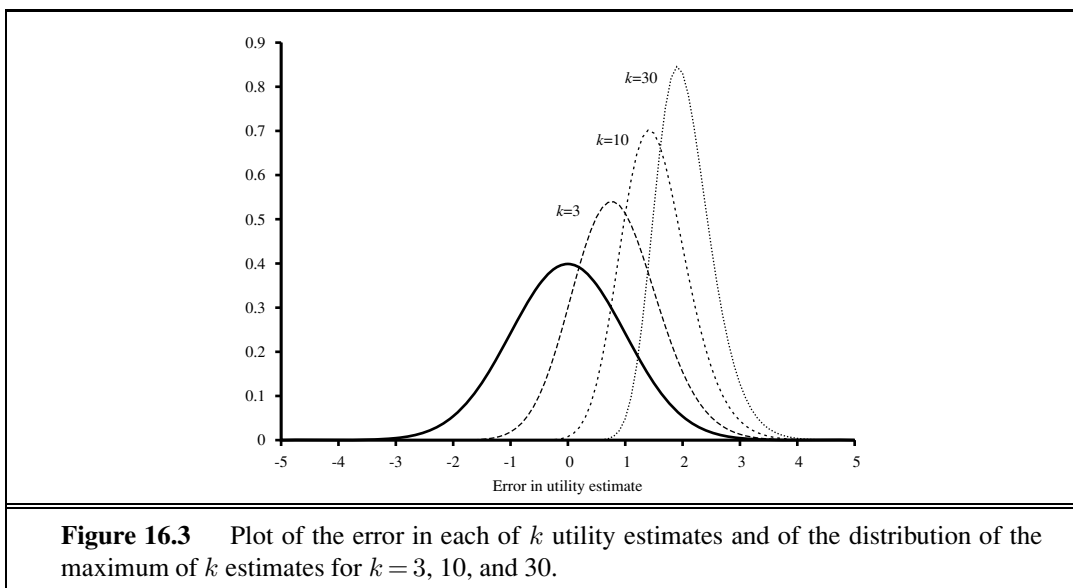
$$a^* = \operatorname{argmax}_a EU(a|\mathbf{e}) .$$

If we have calculated the expected utility correctly according to our probability model, and if the probability model correctly reflects the underlying stochastic processes that generate the outcomes, then, on average, we will get the utility we expect if the whole process is repeated many times.

UNBIASED

In reality, however, our model usually oversimplifies the real situation, either because we don’t know enough (e.g., when making a complex investment decision) or because the computation of the true expected utility is too difficult (e.g., when estimating the utility of successor states of the root node in backgammon). In that case, we are really working with *estimates*  $\widehat{EU}(a|\mathbf{e})$  of the true expected utility. We will assume, kindly perhaps, that the estimates are **unbiased**, that is, the expected value of the error,  $E(\widehat{EU}(a|\mathbf{e}) - EU(a|\mathbf{e}))$ , is zero. In that case, it still seems reasonable to choose the action with the highest estimated utility and to expect to receive that utility, on average, when the action is executed.

Unfortunately, the real outcome will usually be significantly *worse* than we estimated, even though the estimate was unbiased! To see why, consider a decision problem in which there are  $k$  choices, each of which has true estimated utility of 0. Suppose that the error in each utility estimate has zero mean and standard deviation of 1, shown as the bold curve in Figure 16.3. Now, as we actually start to generate the estimates, some of the errors will be negative (pessimistic) and some will be positive (optimistic). Because we select the action with the *highest* utility estimate, we are obviously favoring the overly optimistic estimates, and that is the source of the bias. It is a straightforward matter to calculate the distribution of the maximum of the  $k$  estimates (see Exercise 16.11) and hence quantify the extent of our disappointment. The curve in Figure 16.3 for  $k=3$  has a mean around 0.85, so the average disappointment will be about 85% of the standard deviation in the utility estimates.



With more choices, extremely optimistic estimates are more likely to arise: for  $k = 30$ , the disappointment will be around twice the standard deviation in the estimates.

## OPTIMIZER'S CURSE

This tendency for the estimated expected utility of the best choice to be too high is called the **optimizer's curse** (Smith and Winkler, 2006). It afflicts even the most seasoned decision analysts and statisticians. Serious manifestations include believing that an exciting new drug that has cured 80% patients in a trial will cure 80% of patients (it's been chosen from  $k =$  thousands of candidate drugs) or that a mutual fund advertised as having above-average returns will continue to have them (it's been chosen to appear in the advertisement out of  $k =$  dozens of funds in the company's overall portfolio). It can even be the case that what appears to be the best choice may not be, if the variance in the utility estimate is high: a drug, selected from thousands tried, that has cured 9 of 10 patients is probably *worse* than one that has cured 800 of 1000.

The optimizer's curse crops up everywhere because of the ubiquity of utility-maximizing selection processes, so taking the utility estimates at face value is a bad idea. We can avoid the curse by using an explicit probability model  $\mathbf{P}(\widehat{EU} \mid EU)$  of the error in the utility estimates. Given this model and a prior  $\mathbf{P}(EU)$  on what we might reasonably expect the utilities to be, we treat the utility estimate, once obtained, as evidence and compute the posterior distribution for the true utility using Bayes' rule.

### 16.3.4 Human judgment and irrationality

NORMATIVE THEORY  
DESCRIPTIVE  
THEORY

Decision theory is a **normative theory**: it describes how a rational agent *should* act. A **descriptive theory**, on the other hand, describes how actual agents—for example, humans—really do act. The application of economic theory would be greatly enhanced if the two coincided, but there appears to be some experimental evidence to the contrary. The evidence suggests that humans are “predictably irrational” (Ariely, 2009).

The best-known problem is the Allais paradox (Allais, 1953). People are given a choice between lotteries  $A$  and  $B$  and then between  $C$  and  $D$ , which have the following prizes:

- $A$  : 80% chance of \$4000

$C$  : 20% chance of \$4000
- $B$  : 100% chance of \$3000

$D$  : 25% chance of \$3000

CERTAINTY EFFECT

REGRET

Most people consistently prefer  $B$  over  $A$  (taking the sure thing), and  $C$  over  $D$  (taking the higher EMV). The normative analysis disagrees! We can see this most easily if we use the freedom implied by Equation (16.2) to set  $U(\$0) = 0$ . In that case, then  $B \succ A$  implies that  $U(\$3000) > 0.8U(\$4000)$ , whereas  $C \succ D$  implies exactly the reverse. In other words, there is no utility function that is consistent with these choices. One explanation for the apparently irrational preferences is the **certainty effect** (Kahneman and Tversky, 1979): people are strongly attracted to gains that are certain. There are several reasons why this may be so. First, people may prefer to reduce their computational burden; by choosing certain outcomes, they don't have to compute with probabilities. But the effect persists even when the computations involved are very easy ones. Second, people may distrust the legitimacy of the stated probabilities. I trust that a coin flip is roughly 50/50 if I have control over the coin and the flip, but I may distrust the result if the flip is done by someone with a vested interest in the outcome.<sup>6</sup> In the presence of distrust, it might be better to go for the sure thing.<sup>7</sup> Third, people may be accounting for their emotional state as well as their financial state. People know they would experience **regret** if they gave up a certain reward ( $B$ ) for an 80% chance at a higher reward and then lost. In other words, if  $A$  is chosen, there is a 20% chance of getting no money *and feeling like a complete idiot*, which is worse than just getting no money. So perhaps people who choose  $B$  over  $A$  and  $C$  over  $D$  are not being irrational; they are just saying that they are willing to give up \$200 of EMV to avoid a 20% chance of feeling like an idiot.

A related problem is the Ellsberg paradox. Here the prizes are fixed, but the probabilities are underconstrained. Your payoff will depend on the color of a ball chosen from an urn. You are told that the urn contains 1/3 red balls, and 2/3 either black or yellow balls, but you don't know how many black and how many yellow. Again, you are asked whether you prefer lottery  $A$  or  $B$ ; and then  $C$  or  $D$ :

- $A$  : \$100 for a red ball

$C$  : \$100 for a red or yellow ball
- $B$  : \$100 for a black ball

$D$  : \$100 for a black or yellow ball .

AMBIGUITY  
AVERSION

It should be clear that if you think there are more red than black balls then you should prefer  $A$  over  $B$  and  $C$  over  $D$ ; if you think there are fewer red than black you should prefer the opposite. But it turns out that most people prefer  $A$  over  $B$  and also prefer  $D$  over  $C$ , even though there is no state of the world for which this is rational. It seems that people have **ambiguity aversion**:  $A$  gives you a 1/3 chance of winning, while  $B$  could be anywhere between 0 and 2/3. Similarly,  $D$  gives you a 2/3 chance, while  $C$  could be anywhere between 1/3 and 3/3. Most people elect the known probability rather than the unknown unknowns.

<sup>6</sup> For example, the mathematician/magician Persi Diaconis can make a coin flip come out the way he wants every time (Landhuis, 2004).  
<sup>7</sup> Even the sure thing may not be certain. Despite cast-iron promises, we have not yet received that \$27,000,000 from the Nigerian bank account of a previously unknown deceased relative.



## FRAMING EFFECT

Yet another problem is that the exact wording of a decision problem can have a big impact on the agent's choices; this is called the **framing effect**. Experiments show that people like a medical procedure that it is described as having a "90% survival rate" about twice as much as one described as having a "10% death rate," even though these two statements mean exactly the same thing. This discrepancy in judgment has been found in multiple experiments and is about the same whether the subjects were patients in a clinic, statistically sophisticated business school students, or experienced doctors.

## ANCHORING EFFECT

People feel more comfortable making *relative* utility judgments rather than absolute ones. I may have little idea how much I might enjoy the various wines offered by a restaurant. The restaurant takes advantage of this by offering a \$200 bottle that it knows nobody will buy, but which serves to skew upward the customer's estimate of the value of all wines and make the \$55 bottle seem like a bargain. This is called the **anchoring effect**.

If human informants insist on contradictory preference judgments, there is nothing that automated agents can do to be consistent with them. Fortunately, preference judgments made by humans are often open to revision in the light of further consideration. Paradoxes like the Allais paradox are greatly reduced (but not eliminated) if the choices are explained better. In work at the Harvard Business School on assessing the utility of money, Keeney and Raiffa (1976, p. 210) found the following:

Subjects tend to be too risk-averse in the small and therefore . . . the fitted utility functions exhibit unacceptably large risk premiums for lotteries with a large spread. . . . Most of the subjects, however, can reconcile their inconsistencies and feel that they have learned an important lesson about how they want to behave. As a consequence, some subjects cancel their automobile collision insurance and take out more term insurance on their lives.

EVOLUTIONARY  
PSYCHOLOGY

The evidence for human irrationality is also questioned by researchers in the field of **evolutionary psychology**, who point to the fact that our brain's decision-making mechanisms did not evolve to solve word problems with probabilities and prizes stated as decimal numbers. Let us grant, for the sake of argument, that the brain has built-in neural mechanism for computing with probabilities and utilities, or something functionally equivalent; if so, the required inputs would be obtained through accumulated experience of outcomes and rewards rather than through linguistic presentations of numerical values. It is far from obvious that we can directly access the brain's built-in neural mechanisms by presenting decision problems in linguistic/numerical form. The very fact that different wordings of the *same decision problem* elicit different choices suggests that the decision problem itself is not getting through. Spurred by this observation, psychologists have tried presenting problems in uncertain reasoning and decision making in "evolutionarily appropriate" forms; for example, instead of saying "90% survival rate," the experimenter might show 100 stick-figure animations of the operation, where the patient dies in 10 of them and survives in 90. (Boredom is a complicating factor in these experiments!) With decision problems posed in this way, people seem to be much closer to rational behavior than previously suspected.

## 16.4 MULTIATTRIBUTE UTILITY FUNCTIONS

### MULTIATTRIBUTE UTILITY THEORY

Decision making in the field of public policy involves high stakes, in both money and lives. For example, in deciding what levels of harmful emissions to allow from a power plant, policy makers must weigh the prevention of death and disability against the benefit of the power and the economic burden of mitigating the emissions. Siting a new airport requires consideration of the disruption caused by construction; the cost of land; the distance from centers of population; the noise of flight operations; safety issues arising from local topography and weather conditions; and so on. Problems like these, in which outcomes are characterized by two or more attributes, are handled by **multiattribute utility theory**.

We will call the attributes  $\mathbf{X} = X_1, \dots, X_n$ ; a complete vector of assignments will be  $\mathbf{x} = \langle x_1, \dots, x_n \rangle$ , where each  $x_i$  is either a numeric value or a discrete value with an assumed ordering on values. We will assume that higher values of an attribute correspond to higher utilities, all other things being equal. For example, if we choose *AbsenceOfNoise* as an attribute in the airport problem, then the greater its value, the better the solution.<sup>8</sup> We begin by examining cases in which decisions can be made *without* combining the attribute values into a single utility value. Then we look at cases in which the utilities of attribute combinations can be specified very concisely.

### 16.4.1 Dominance

#### STRICT DOMINANCE

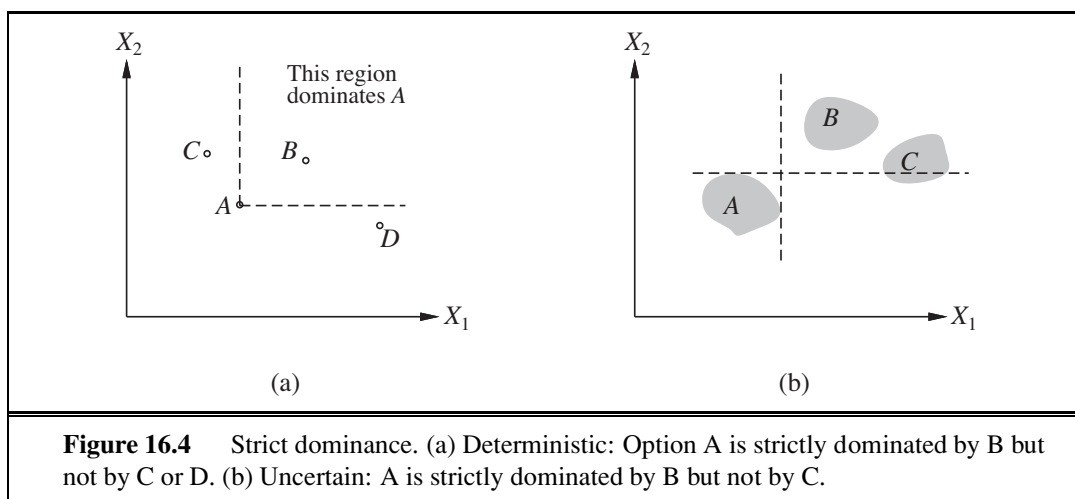
Suppose that airport site  $S_1$  costs less, generates less noise pollution, and is safer than site  $S_2$ . One would not hesitate to reject  $S_2$ . We then say that there is **strict dominance** of  $S_1$  over  $S_2$ . In general, if an option is of lower value on all attributes than some other option, it need not be considered further. Strict dominance is often very useful in narrowing down the field of choices to the real contenders, although it seldom yields a unique choice. Figure 16.4(a) shows a schematic diagram for the two-attribute case.

That is fine for the deterministic case, in which the attribute values are known for sure. What about the general case, where the outcomes are uncertain? A direct analog of strict dominance can be constructed, where, despite the uncertainty, all possible concrete outcomes for  $S_1$  strictly dominate all possible outcomes for  $S_2$ . (See Figure 16.4(b).) Of course, this will probably occur even less often than in the deterministic case.

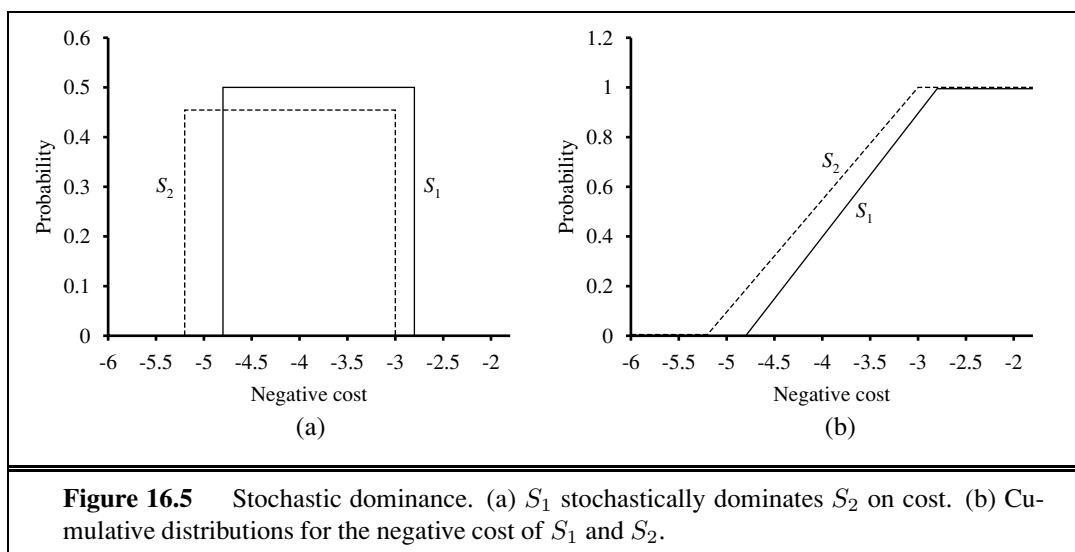
#### STOCHASTIC DOMINANCE

Fortunately, there is a more useful generalization called **stochastic dominance**, which occurs very frequently in real problems. Stochastic dominance is easiest to understand in the context of a single attribute. Suppose we believe that the cost of siting the airport at  $S_1$  is uniformly distributed between \$2.8 billion and \$4.8 billion and that the cost at  $S_2$  is uniformly distributed between \$3 billion and \$5.2 billion. Figure 16.5(a) shows these distributions, with cost plotted as a negative value. Then, given only the information that utility decreases with

<sup>8</sup> In some cases, it may be necessary to subdivide the range of values so that utility varies monotonically within each range. For example, if the *RoomTemperature* attribute has a utility peak at 70°F, we would split it into two attributes measuring the difference from the ideal, one colder and one hotter. Utility would then be monotonically increasing in each attribute.



**Figure 16.4** Strict dominance. (a) Deterministic: Option A is strictly dominated by B but not by C or D. (b) Uncertain: A is strictly dominated by B but not by C.



**Figure 16.5** Stochastic dominance. (a)  $S_1$  stochastically dominates  $S_2$  on cost. (b) Cumulative distributions for the negative cost of  $S_1$  and  $S_2$ .

cost, we can say that  $S_1$  stochastically dominates  $S_2$  (i.e.,  $S_2$  can be discarded). It is important to note that this does *not* follow from comparing the expected costs. For example, if we knew the cost of  $S_1$  to be *exactly* \$3.8 billion, then we would be *unable* to make a decision without additional information on the utility of money. (It might seem odd that *more* information on the cost of  $S_1$  could make the agent *less* able to decide. The paradox is resolved by noting that in the absence of exact cost information, the decision is easier to make but is more likely to be wrong.)

The exact relationship between the attribute distributions needed to establish stochastic dominance is best seen by examining the **cumulative distributions**, shown in Figure 16.5(b). (See also Appendix A.) The cumulative distribution measures the probability that the cost is less than or equal to any given amount—that is, it integrates the original distribution. If the cumulative distribution for  $S_1$  is always to the right of the cumulative distribution for  $S_2$ ,

then, stochastically speaking,  $S_1$  is cheaper than  $S_2$ . Formally, if two actions  $A_1$  and  $A_2$  lead to probability distributions  $p_1(x)$  and  $p_2(x)$  on attribute  $X$ , then  $A_1$  stochastically dominates  $A_2$  on  $X$  if

$$\forall x \int_{-\infty}^x p_1(x') dx' \leq \int_{-\infty}^x p_2(x') dx' .$$



The relevance of this definition to the selection of optimal decisions comes from the following property: *if  $A_1$  stochastically dominates  $A_2$ , then for any monotonically nondecreasing utility function  $U(x)$ , the expected utility of  $A_1$  is at least as high as the expected utility of  $A_2$ .* Hence, if an action is stochastically dominated by another action on all attributes, then it can be discarded.

The stochastic dominance condition might seem rather technical and perhaps not so easy to evaluate without extensive probability calculations. In fact, it can be decided very easily in many cases. Suppose, for example, that the construction transportation cost depends on the distance to the supplier. The cost itself is uncertain, but the greater the distance, the greater the cost. If  $S_1$  is closer than  $S_2$ , then  $S_1$  will dominate  $S_2$  on cost. Although we will not present them here, there exist algorithms for propagating this kind of qualitative information among uncertain variables in **qualitative probabilistic networks**, enabling a system to make rational decisions based on stochastic dominance, without using any numeric values.

QUALITATIVE  
PROBABILISTIC  
NETWORKS

### 16.4.2 Preference structure and multiattribute utility

Suppose we have  $n$  attributes, each of which has  $d$  distinct possible values. To specify the complete utility function  $U(x_1, \dots, x_n)$ , we need  $d^n$  values in the worst case. Now, the worst case corresponds to a situation in which the agent's preferences have no regularity at all. Multiattribute utility theory is based on the supposition that the preferences of typical agents have much more structure than that. The basic approach is to identify regularities in the preference behavior we would expect to see and to use what are called **representation theorems** to show that an agent with a certain kind of preference structure has a utility function

$$U(x_1, \dots, x_n) = F[f_1(x_1), \dots, f_n(x_n)] ,$$

where  $F$  is, we hope, a simple function such as addition. Notice the similarity to the use of Bayesian networks to decompose the joint probability of several random variables.

REPRESENTATION  
THEOREM

#### Preferences without uncertainty

Let us begin with the deterministic case. Remember that for deterministic environments the agent has a value function  $V(x_1, \dots, x_n)$ ; the aim is to represent this function concisely. The basic regularity that arises in deterministic preference structures is called **preference independence**. Two attributes  $X_1$  and  $X_2$  are preferentially independent of a third attribute  $X_3$  if the preference between outcomes  $\langle x_1, x_2, x_3 \rangle$  and  $\langle x'_1, x'_2, x_3 \rangle$  does not depend on the particular value  $x_3$  for attribute  $X_3$ .

PREFERENCE  
INDEPENDENCE

Going back to the airport example, where we have (among other attributes) *Noise*, *Cost*, and *Deaths* to consider, one may propose that *Noise* and *Cost* are preferentially inde-

MUTUAL  
PREFERENTIAL  
INDEPENDENCE

pendent of *Deaths*. For example, if we prefer a state with 20,000 people residing in the flight path and a construction cost of \$4 billion over a state with 70,000 people residing in the flight path and a cost of \$3.7 billion when the safety level is 0.06 deaths per million passenger miles in both cases, then we would have the same preference when the safety level is 0.12 or 0.03; and the same independence would hold for preferences between any other pair of values for *Noise* and *Cost*. It is also apparent that *Cost* and *Deaths* are preferentially independent of *Noise* and that *Noise* and *Deaths* are preferentially independent of *Cost*. We say that the set of attributes  $\{\textit{Noise}, \textit{Cost}, \textit{Deaths}\}$  exhibits **mutual preferential independence** (MPI). MPI says that, whereas each attribute may be important, it does not affect the way in which one trades off the other attributes against each other.

Mutual preferential independence is something of a mouthful, but thanks to a remarkable theorem due to the economist Gérard Debreu (1960), we can derive from it a very simple form for the agent's value function: *If attributes  $X_1, \dots, X_n$  are mutually preferentially independent, then the agent's preference behavior can be described as maximizing the function*

$$V(x_1, \dots, x_n) = \sum_i V_i(x_i),$$

where each  $V_i$  is a value function referring only to the attribute  $X_i$ . For example, it might well be the case that the airport decision can be made using a value function

$$V(\textit{noise}, \textit{cost}, \textit{deaths}) = -\textit{noise} \times 10^4 - \textit{cost} - \textit{deaths} \times 10^{12}.$$

ADDITIVE VALUE  
FUNCTION

A value function of this type is called an **additive value function**. Additive functions are an extremely natural way to describe an agent's preferences and are valid in many real-world situations. For  $n$  attributes, assessing an additive value function requires assessing  $n$  separate one-dimensional value functions rather than one  $n$ -dimensional function; typically, this represents an exponential reduction in the number of preference experiments that are needed. Even when MPI does not strictly hold, as might be the case at extreme values of the attributes, an additive value function might still provide a good approximation to the agent's preferences. This is especially true when the violations of MPI occur in portions of the attribute ranges that are unlikely to occur in practice.

To understand MPI better, it helps to look at cases where it *doesn't* hold. Suppose you are at a medieval market, considering the purchase of some hunting dogs, some chickens, and some wicker cages for the chickens. The hunting dogs are very valuable, but if you don't have enough cages for the chickens, the dogs will eat the chickens; hence, the tradeoff between dogs and chickens depends strongly on the number of cages, and MPI is violated. The existence of these kinds of interactions among various attributes makes it much harder to assess the overall value function.

### Preferences with uncertainty

When uncertainty is present in the domain, we also need to consider the structure of preferences between lotteries and to understand the resulting properties of utility functions, rather than just value functions. The mathematics of this problem can become quite complicated, so we present just one of the main results to give a flavor of what can be done. The reader is referred to Keeney and Raiffa (1976) for a thorough survey of the field.

UTILITY  
INDEPENDENCE

The basic notion of **utility independence** extends preference independence to cover lotteries: a set of attributes  $\mathbf{X}$  is utility independent of a set of attributes  $\mathbf{Y}$  if preferences between lotteries on the attributes in  $\mathbf{X}$  are independent of the particular values of the attributes in  $\mathbf{Y}$ . A set of attributes is **mutually utility independent** (MUI) if each of its subsets is utility-independent of the remaining attributes. Again, it seems reasonable to propose that the airport attributes are MUI.

MUTUALLY UTILITY  
INDEPENDENTMULTIPLICATIVE  
UTILITY FUNCTION

MUI implies that the agent's behavior can be described using a **multiplicative utility function** (Keeney, 1974). The general form of a multiplicative utility function is best seen by looking at the case for three attributes. For conciseness, we use  $U_i$  to mean  $U_i(x_i)$ :

$$U = k_1U_1 + k_2U_2 + k_3U_3 + k_1k_2U_1U_2 + k_2k_3U_2U_3 + k_3k_1U_3U_1 + k_1k_2k_3U_1U_2U_3.$$

Although this does not look very simple, it contains just three single-attribute utility functions and three constants. In general, an  $n$ -attribute problem exhibiting MUI can be modeled using  $n$  single-attribute utilities and  $n$  constants. Each of the single-attribute utility functions can be developed independently of the other attributes, and this combination will be guaranteed to generate the correct overall preferences. Additional assumptions are required to obtain a purely additive utility function.

## 16.5 DECISION NETWORKS

INFLUENCE DIAGRAM  
DECISION NETWORK

In this section, we look at a general mechanism for making rational decisions. The notation is often called an **influence diagram** (Howard and Matheson, 1984), but we will use the more descriptive term **decision network**. Decision networks combine Bayesian networks with additional node types for actions and utilities. We use airport siting as an example.

### 16.5.1 Representing a decision problem with a decision network

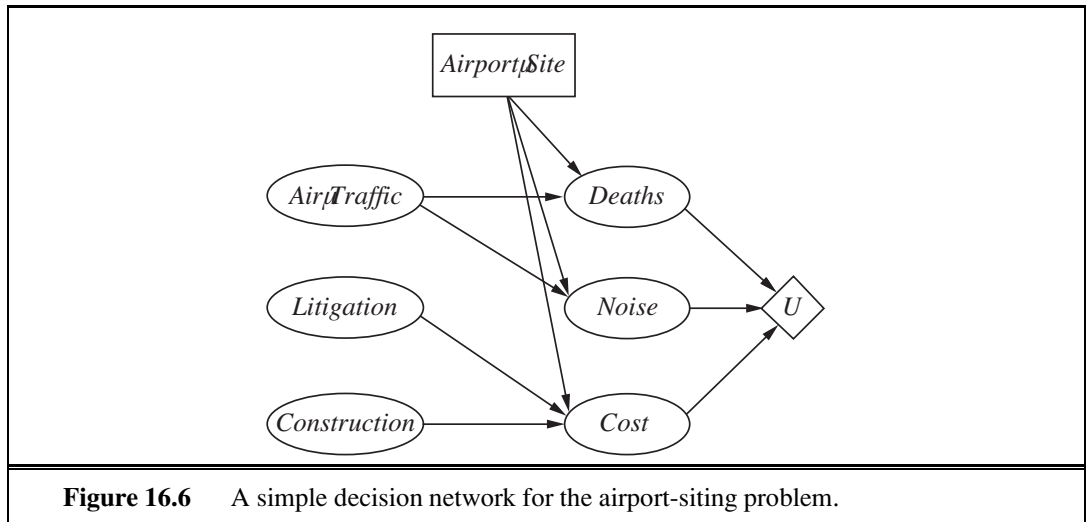
In its most general form, a decision network represents information about the agent's current state, its possible actions, the state that will result from the agent's action, and the utility of that state. It therefore provides a substrate for implementing utility-based agents of the type first introduced in Section 2.4.5. Figure 16.6 shows a decision network for the airport siting problem. It illustrates the three types of nodes used:

CHANCE NODES

- **Chance nodes** (ovals) represent random variables, just as they do in Bayesian networks. The agent could be uncertain about the construction cost, the level of air traffic and the potential for litigation, and the *Deaths*, *Noise*, and total *Cost* variables, each of which also depends on the site chosen. Each chance node has associated with it a conditional distribution that is indexed by the state of the parent nodes. In decision networks, the parent nodes can include decision nodes as well as chance nodes. Note that each of the current-state chance nodes could be part of a large Bayesian network for assessing construction costs, air traffic levels, or litigation potentials.

DECISION NODES

- **Decision nodes** (rectangles) represent points where the decision maker has a choice of



actions. In this case, the *AirportSite* action can take on a different value for each site under consideration. The choice influences the cost, safety, and noise that will result. In this chapter, we assume that we are dealing with a single decision node. Chapter 17 deals with cases in which more than one decision must be made.

## UTILITY NODES

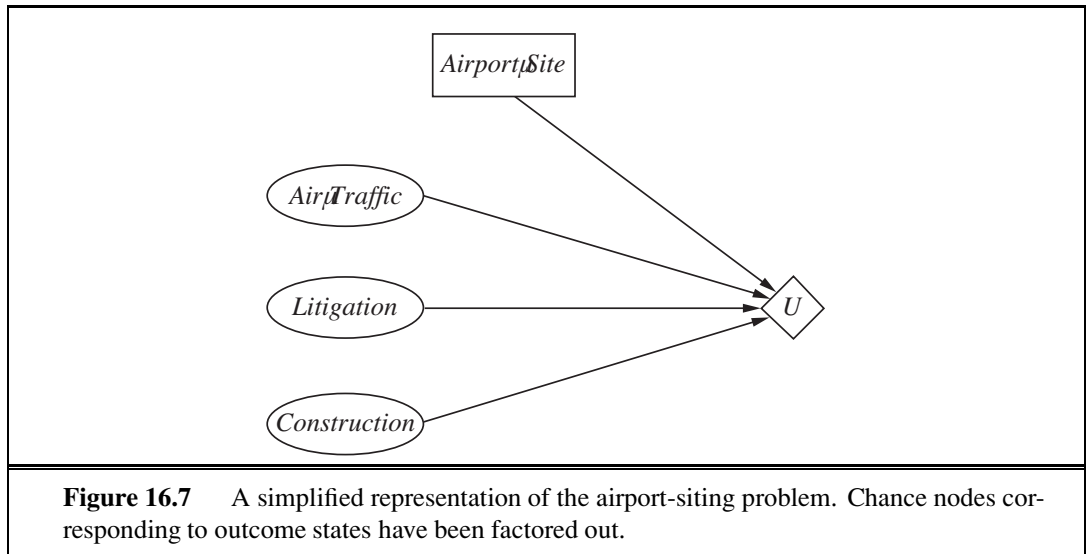
- **Utility nodes** (diamonds) represent the agent's utility function.<sup>9</sup> The utility node has as parents all variables describing the outcome that directly affect utility. Associated with the utility node is a description of the agent's utility as a function of the parent attributes. The description could be just a tabulation of the function, or it might be a parameterized additive or linear function of the attribute values.

A simplified form is also used in many cases. The notation remains identical, but the chance nodes describing the outcome state are omitted. Instead, the utility node is connected directly to the current-state nodes and the decision node. In this case, rather than representing a utility function on outcome states, the utility node represents the *expected* utility associated with each action, as defined in Equation (16.1) on page 611; that is, the node is associated with an **action-utility function** (also known as a **Q-function** in reinforcement learning, as described in Chapter 21). Figure 16.7 shows the action-utility representation of the airport siting problem.

## ACTION-UTILITY FUNCTION

Notice that, because the *Noise*, *Deaths*, and *Cost* chance nodes in Figure 16.6 refer to future states, they can never have their values set as evidence variables. Thus, the simplified version that omits these nodes can be used whenever the more general form can be used. Although the simplified form contains fewer nodes, the omission of an explicit description of the outcome of the siting decision means that it is less flexible with respect to changes in circumstances. For example, in Figure 16.6, a change in aircraft noise levels can be reflected by a change in the conditional probability table associated with the *Noise* node, whereas a change in the weight accorded to noise pollution in the utility function can be reflected by

<sup>9</sup> These nodes are also called **value nodes** in the literature.



a change in the utility table. In the action-utility diagram, Figure 16.7, on the other hand, all such changes have to be reflected by changes to the action-utility table. Essentially, the action-utility formulation is a *compiled* version of the original formulation.

### 16.5.2 Evaluating decision networks

Actions are selected by evaluating the decision network for each possible setting of the decision node. Once the decision node is set, it behaves exactly like a chance node that has been set as an evidence variable. The algorithm for evaluating decision networks is the following:

1. Set the evidence variables for the current state.
2. For each possible value of the decision node:
  - (a) Set the decision node to that value.
  - (b) Calculate the posterior probabilities for the parent nodes of the utility node, using a standard probabilistic inference algorithm.
  - (c) Calculate the resulting utility for the action.
3. Return the action with the highest utility.

This is a straightforward extension of the Bayesian network algorithm and can be incorporated directly into the agent design given in Figure 13.1 on page 484. We will see in Chapter 17 that the possibility of executing several actions in sequence makes the problem much more interesting.

## 16.6 THE VALUE OF INFORMATION

In the preceding analysis, we have assumed that all relevant information, or at least all available information, is provided to the agent before it makes its decision. In practice, this is





#### INFORMATION VALUE THEORY

hardly ever the case. *One of the most important parts of decision making is knowing what questions to ask.* For example, a doctor cannot expect to be provided with the results of *all possible* diagnostic tests and questions at the time a patient first enters the consulting room.<sup>10</sup> Tests are often expensive and sometimes hazardous (both directly and because of associated delays). Their importance depends on two factors: whether the test results would lead to a significantly better treatment plan, and how likely the various test results are.

This section describes **information value theory**, which enables an agent to choose what information to acquire. We assume that, prior to selecting a “real” action represented by the decision node, the agent can acquire the value of any of the potentially observable chance variables in the model. Thus, information value theory involves a simplified form of sequential decision making—simplified because the observation actions affect only the agent’s **belief state**, not the external physical state. The value of any particular observation must derive from the potential to affect the agent’s eventual physical action; and this potential can be estimated directly from the decision model itself.

### 16.6.1 A simple example

Suppose an oil company is hoping to buy one of  $n$  indistinguishable blocks of ocean-drilling rights. Let us assume further that exactly one of the blocks contains oil worth  $C$  dollars, while the others are worthless. The asking price of each block is  $C/n$  dollars. If the company is risk-neutral, then it will be indifferent between buying a block and not buying one.

Now suppose that a seismologist offers the company the results of a survey of block number 3, which indicates definitively whether the block contains oil. How much should the company be willing to pay for the information? The way to answer this question is to examine what the company would do if it had the information:

- With probability  $1/n$ , the survey will indicate oil in block 3. In this case, the company will buy block 3 for  $C/n$  dollars and make a profit of  $C - C/n = (n - 1)C/n$  dollars.
- With probability  $(n - 1)/n$ , the survey will show that the block contains no oil, in which case the company will buy a different block. Now the probability of finding oil in one of the other blocks changes from  $1/n$  to  $1/(n - 1)$ , so the company makes an expected profit of  $C/(n - 1) - C/n = C/n(n - 1)$  dollars.

Now we can calculate the expected profit, given the survey information:

$$\frac{1}{n} \times \frac{(n - 1)C}{n} + \frac{n - 1}{n} \times \frac{C}{n(n - 1)} = C/n.$$

Therefore, the company should be willing to pay the seismologist up to  $C/n$  dollars for the information: the information is worth as much as the block itself.

The value of information derives from the fact that *with* the information, one’s course of action can be changed to suit the *actual* situation. One can discriminate according to the situation, whereas without the information, one has to do what’s best on average over the possible situations. In general, the value of a given piece of information is defined to be the difference in expected value between best actions before and after information is obtained.

<sup>10</sup> In the United States, the only question that is always asked beforehand is whether the patient has insurance.

### 16.6.2 A general formula for perfect information

It is simple to derive a general mathematical formula for the value of information. We assume that exact evidence can be obtained about the value of some random variable  $E_j$  (that is, we learn  $E_j = e_j$ ), so the phrase **value of perfect information** (VPI) is used.<sup>11</sup>

Let the agent's initial evidence be  $\mathbf{e}$ . Then the value of the current best action  $\alpha$  is defined by

$$EU(\alpha|\mathbf{e}) = \max_a \sum_{s'} P(\text{RESULT}(a) = s' | a, \mathbf{e}) U(s') ,$$

and the value of the new best action (after the new evidence  $E_j = e_j$  is obtained) will be

$$EU(\alpha_{e_j}|\mathbf{e}, e_j) = \max_a \sum_{s'} P(\text{RESULT}(a) = s' | a, \mathbf{e}, e_j) U(s') .$$

But  $E_j$  is a random variable whose value is *currently* unknown, so to determine the value of discovering  $E_j$ , given current information  $\mathbf{e}$  we must average over all possible values  $e_{jk}$  that we might discover for  $E_j$ , using our *current* beliefs about its value:

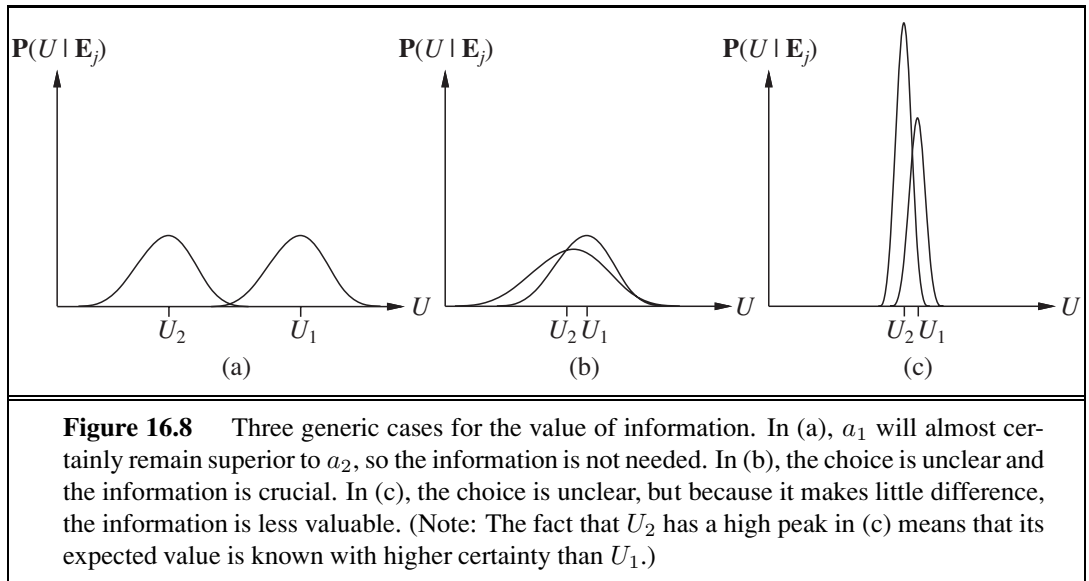
$$VPI_{\mathbf{e}}(E_j) = \left( \sum_k P(E_j = e_{jk}|\mathbf{e}) EU(\alpha_{e_{jk}}|\mathbf{e}, E_j = e_{jk}) \right) - EU(\alpha|\mathbf{e}) .$$

To get some intuition for this formula, consider the simple case where there are only two actions,  $a_1$  and  $a_2$ , from which to choose. Their current expected utilities are  $U_1$  and  $U_2$ . The information  $E_j = e_{jk}$  will yield some new expected utilities  $U'_1$  and  $U'_2$  for the actions, but before we obtain  $E_j$ , we will have some probability distributions over the possible values of  $U'_1$  and  $U'_2$  (which we assume are independent).

Suppose that  $a_1$  and  $a_2$  represent two different routes through a mountain range in winter.  $a_1$  is a nice, straight highway through a low pass, and  $a_2$  is a winding dirt road over the top. Just given this information,  $a_1$  is clearly preferable, because it is quite possible that  $a_2$  is blocked by avalanches, whereas it is unlikely that anything blocks  $a_1$ .  $U_1$  is therefore clearly higher than  $U_2$ . It is possible to obtain satellite reports  $E_j$  on the actual state of each road that would give new expectations,  $U'_1$  and  $U'_2$ , for the two crossings. The distributions for these expectations are shown in Figure 16.8(a). Obviously, in this case, it is not worth the expense of obtaining satellite reports, because it is unlikely that the information derived from them will change the plan. With no change, information has no value.

Now suppose that we are choosing between two different winding dirt roads of slightly different lengths and we are carrying a seriously injured passenger. Then, even when  $U_1$  and  $U_2$  are quite close, the distributions of  $U'_1$  and  $U'_2$  are very broad. There is a significant possibility that the second route will turn out to be clear while the first is blocked, and in this

<sup>11</sup> There is no loss of expressiveness in requiring perfect information. Suppose we wanted to model the case in which we become somewhat more certain about a variable. We can do that by introducing *another* variable about which we learn perfect information. For example, suppose we initially have broad uncertainty about the variable *Temperature*. Then we gain the perfect knowledge *Thermometer* = 37; this gives us imperfect information about the true *Temperature*, and the uncertainty due to measurement error is encoded in the sensor model  $\mathbf{P}(\text{Thermometer} | \text{Temperature})$ . See Exercise 16.17 for another example.



case the difference in utilities will be very high. The VPI formula indicates that it might be worthwhile getting the satellite reports. Such a situation is shown in Figure 16.8(b).

Finally, suppose that we are choosing between the two dirt roads in summertime, when blockage by avalanches is unlikely. In this case, satellite reports might show one route to be more scenic than the other because of flowering alpine meadows, or perhaps wetter because of errant streams. It is therefore quite likely that we would change our plan if we had the information. In this case, however, the difference in value between the two routes is still likely to be very small, so we will not bother to obtain the reports. This situation is shown in Figure 16.8(c).



In sum, *information has value to the extent that it is likely to cause a change of plan and to the extent that the new plan will be significantly better than the old plan.*

### 16.6.3 Properties of the value of information

One might ask whether it is possible for information to be deleterious: can it actually have negative expected value? Intuitively, one should expect this to be impossible. After all, one could in the worst case just ignore the information and pretend that one has never received it. This is confirmed by the following theorem, which applies to any decision-theoretic agent:



*The expected value of information is nonnegative:*

$$\forall \mathbf{e}, E_j \quad \text{VPI}_{\mathbf{e}}(E_j) \geq 0.$$

The theorem follows directly from the definition of VPI, and we leave the proof as an exercise (Exercise 16.18). It is, of course, a theorem about *expected* value, not *actual* value. Additional information can easily lead to a plan that *turns out to be* worse than the original plan if the information happens to be misleading. For example, a medical test that gives a false positive result may lead to unnecessary surgery; but that does not mean that the test shouldn't be done.

It is important to remember that VPI depends on the current state of information, which is why it is subscripted. It can change as more information is acquired. For any given piece of evidence  $E_j$ , the value of acquiring it can go down (e.g., if another variable strongly constrains the posterior for  $E_j$ ) or up (e.g., if another variable provides a clue on which  $E_j$  builds, enabling a new and better plan to be devised). Thus, VPI is not additive. That is,

$$VPI_e(E_j, E_k) \neq VPI_e(E_j) + VPI_e(E_k) \quad (\text{in general}) .$$

VPI is, however, order independent. That is,

$$VPI_e(E_j, E_k) = VPI_e(E_j) + VPI_{e,e_j}(E_k) = VPI_e(E_k) + VPI_{e,e_k}(E_j) .$$

Order independence distinguishes sensing actions from ordinary actions and simplifies the problem of calculating the value of a sequence of sensing actions.

### 16.6.4 Implementation of an information-gathering agent

A sensible agent should ask questions in a reasonable order, should avoid asking questions that are irrelevant, should take into account the importance of each piece of information in relation to its cost, and should stop asking questions when that is appropriate. All of these capabilities can be achieved by using the value of information as a guide.

Figure 16.9 shows the overall design of an agent that can gather information intelligently before acting. For now, we assume that with each observable evidence variable  $E_j$ , there is an associated cost,  $Cost(E_j)$ , which reflects the cost of obtaining the evidence through tests, consultants, questions, or whatever. The agent requests what appears to be the most efficient observation in terms of utility gain per unit cost. We assume that the result of the action  $Request(E_j)$  is that the next percept provides the value of  $E_j$ . If no observation is worth its cost, the agent selects a “real” action.

The agent algorithm we have described implements a form of information gathering that is called **myopic**. This is because it uses the VPI formula shortsightedly, calculating the value of information as if only a single evidence variable will be acquired. Myopic control is based on the same heuristic idea as greedy search and often works well in practice. (For example, it has been shown to outperform expert physicians in selecting diagnostic tests.)

MYOPIC

**function** INFORMATION-GATHERING-AGENT(*percept*) **returns** an *action*

**persistent:**  $D$ , a decision network

integrate *percept* into  $D$

$j \leftarrow$  the value that maximizes  $VPI(E_j) / Cost(E_j)$

**if**  $VPI(E_j) > Cost(E_j)$

**return** REQUEST( $E_j$ )

**else return** the best action from  $D$

**Figure 16.9** Design of a simple information-gathering agent. The agent works by repeatedly selecting the observation with the highest information value, until the cost of the next observation is greater than its expected benefit.

However, if there is no single evidence variable that will help a lot, a myopic agent might hastily take an action when it would have been better to request two or more variables first and then take action. A better approach in this situation would be to construct a *conditional plan* (as described in Section 11.3.2) that asks for variable values and takes different next steps depending on the answer.

One final consideration is the effect a series of questions will have on a human respondent. People may respond better to a series of questions if they “make sense,” so some expert systems are built to take this into account, asking questions in an order that maximizes the total utility of the system and human rather than an order that maximizes value of information.

## 16.7 DECISION-THEORETIC EXPERT SYSTEMS

---

DECISION ANALYSIS    The field of **decision analysis**, which evolved in the 1950s and 1960s, studies the application of decision theory to actual decision problems. It is used to help make rational decisions in important domains where the stakes are high, such as business, government, law, military strategy, medical diagnosis and public health, engineering design, and resource management. The process involves a careful study of the possible actions and outcomes, as well as the preferences placed on each outcome. It is traditional in decision analysis to talk about two roles: the **decision maker** states preferences between outcomes, and the **decision analyst** enumerates the possible actions and outcomes and elicits preferences from the decision maker to determine the best course of action. Until the early 1980s, the main purpose of decision analysis was to help humans make decisions that actually reflect their own preferences. As more and more decision processes become automated, decision analysis is increasingly used to ensure that the automated processes are behaving as desired.

DECISION MAKER

DECISION ANALYST

Early expert system research concentrated on answering questions, rather than on making decisions. Those systems that did recommend actions rather than providing opinions on matters of fact generally did so using condition-action rules, rather than with explicit representations of outcomes and preferences. The emergence of Bayesian networks in the late 1980s made it possible to build large-scale systems that generated sound probabilistic inferences from evidence. The addition of decision networks means that expert systems can be developed that recommend optimal decisions, reflecting the preferences of the agent as well as the available evidence.

A system that incorporates utilities can avoid one of the most common pitfalls associated with the consultation process: confusing likelihood and importance. A common strategy in early medical expert systems, for example, was to rank possible diagnoses in order of likelihood and report the most likely. Unfortunately, this can be disastrous! For the majority of patients in general practice, the two most *likely* diagnoses are usually “There’s nothing wrong with you” and “You have a bad cold,” but if the third most likely diagnosis for a given patient is lung cancer, that’s a serious matter. Obviously, a testing or treatment plan should depend both on probabilities and utilities. Current medical expert systems can take into account the value of information to recommend tests, and then describe a differential diagnosis.

We now describe the knowledge engineering process for decision-theoretic expert systems. As an example we consider the problem of selecting a medical treatment for a kind of congenital heart disease in children (see Lucas, 1996).

AORTIC  
COARCTATION

About 0.8% of children are born with a heart anomaly, the most common being **aortic coarctation** (a constriction of the aorta). It can be treated with surgery, angioplasty (expanding the aorta with a balloon placed inside the artery), or medication. The problem is to decide what treatment to use and when to do it: the younger the infant, the greater the risks of certain treatments, but one mustn't wait too long. A decision-theoretic expert system for this problem can be created by a team consisting of at least one domain expert (a pediatric cardiologist) and one knowledge engineer. The process can be broken down into the following steps:

**Create a causal model.** Determine the possible symptoms, disorders, treatments, and outcomes. Then draw arcs between them, indicating what disorders cause what symptoms, and what treatments alleviate what disorders. Some of this will be well known to the domain expert, and some will come from the literature. Often the model will match well with the informal graphical descriptions given in medical textbooks.

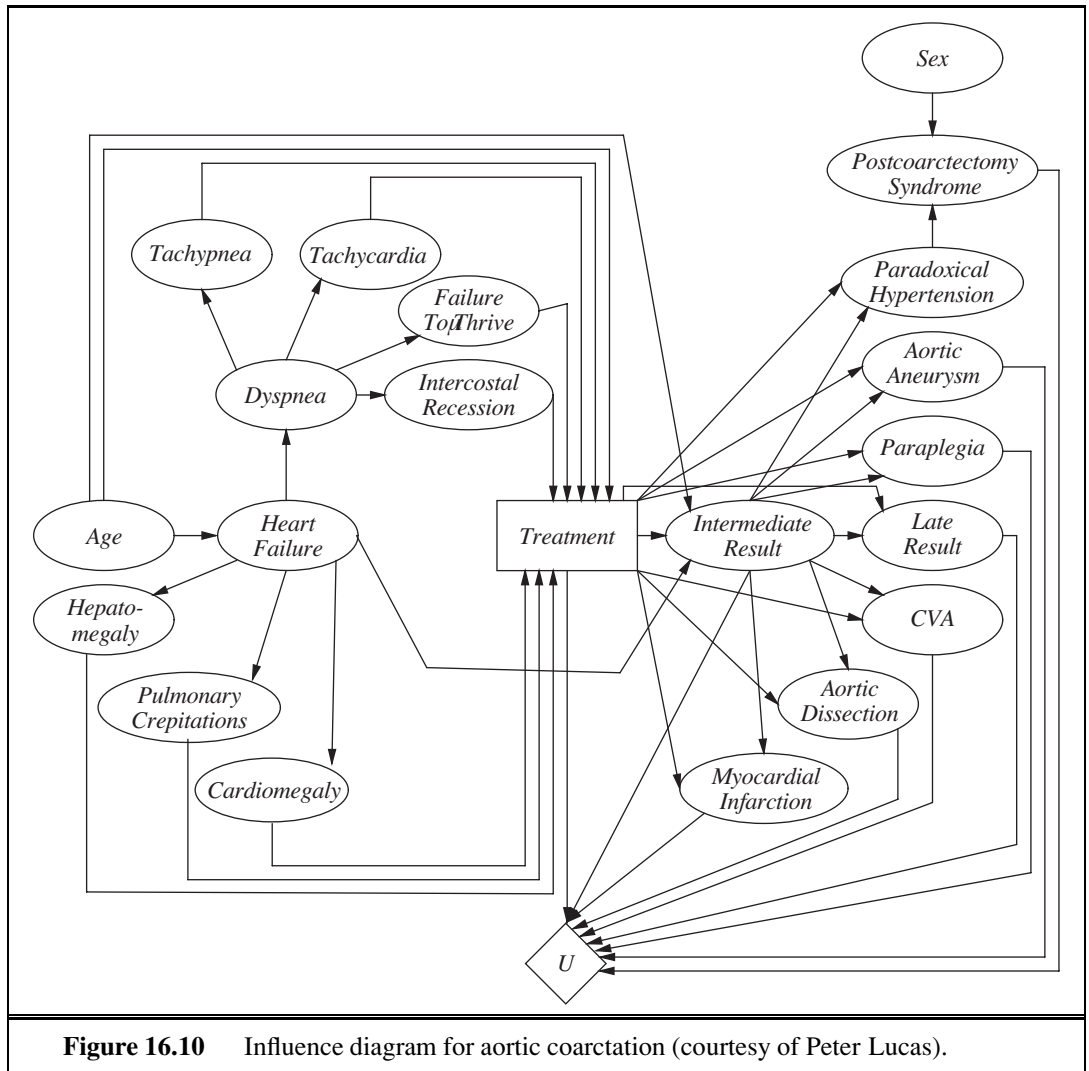
**Simplify to a qualitative decision model.** Since we are using the model to make treatment decisions and not for other purposes (such as determining the joint probability of certain symptom/disorder combinations), we can often simplify by removing variables that are not involved in treatment decisions. Sometimes variables will have to be split or joined to match the expert's intuitions. For example, the original aortic coarctation model had a *Treatment* variable with values *surgery*, *angioplasty*, and *medication*, and a separate variable for *Timing* of the treatment. But the expert had a hard time thinking of these separately, so they were combined, with *Treatment* taking on values such as *surgery in 1 month*. This gives us the model of Figure 16.10.

**Assign probabilities.** Probabilities can come from patient databases, literature studies, or the expert's subjective assessments. Note that a diagnostic system will reason from symptoms and other observations to the disease or other cause of the problems. Thus, in the early years of building these systems, experts were asked for the probability of a cause given an effect. In general they found this difficult to do, and were better able to assess the probability of an effect given a cause. So modern systems usually assess causal knowledge and encode it directly in the Bayesian network structure of the model, leaving the diagnostic reasoning to the Bayesian network inference algorithms (Shachter and Heckerman, 1987).

**Assign utilities.** When there are a small number of possible outcomes, they can be enumerated and evaluated individually using the methods of Section 16.3.1. We would create a scale from best to worst outcome and give each a numeric value, for example 0 for death and 1 for complete recovery. We would then place the other outcomes on this scale. This can be done by the expert, but it is better if the patient (or in the case of infants, the patient's parents) can be involved, because different people have different preferences. If there are exponentially many outcomes, we need some way to combine them using multiattribute utility functions. For example, we may say that the costs of various complications are additive.

**Verify and refine the model.** To evaluate the system we need a set of correct (input, output) pairs; a so-called **gold standard** to compare against. For medical expert systems this usually means assembling the best available doctors, presenting them with a few cases,

GOLD STANDARD



and asking them for their diagnosis and recommended treatment plan. We then see how well the system matches their recommendations. If it does poorly, we try to isolate the parts that are going wrong and fix them. It can be useful to run the system “backward.” Instead of presenting the system with symptoms and asking for a diagnosis, we can present it with a diagnosis such as “heart failure,” examine the predicted probability of symptoms such as tachycardia, and compare with the medical literature.

#### SENSITIVITY ANALYSIS

**Perform sensitivity analysis.** This important step checks whether the best decision is sensitive to small changes in the assigned probabilities and utilities by systematically varying those parameters and running the evaluation again. If small changes lead to significantly different decisions, then it could be worthwhile to spend more resources to collect better data. If all variations lead to the same decision, then the agent will have more confidence that it is the right decision. Sensitivity analysis is particularly important, because one of the main

criticisms of probabilistic approaches to expert systems is that it is too difficult to assess the numerical probabilities required. Sensitivity analysis often reveals that many of the numbers need be specified only very approximately. For example, we might be uncertain about the conditional probability  $P(\text{tachycardia} \mid \text{dyspnea})$ , but if the optimal decision is reasonably robust to small variations in the probability, then our ignorance is less of a concern.

## 16.8 SUMMARY

---

This chapter shows how to combine utility theory with probability to enable an agent to select actions that will maximize its expected performance.

- **Probability theory** describes what an agent should believe on the basis of evidence, **utility theory** describes what an agent wants, and **decision theory** puts the two together to describe what an agent should do.
- We can use decision theory to build a system that makes decisions by considering all possible actions and choosing the one that leads to the best expected outcome. Such a system is known as a **rational agent**.
- Utility theory shows that an agent whose preferences between lotteries are consistent with a set of simple axioms can be described as possessing a utility function; furthermore, the agent selects actions as if maximizing its expected utility.
- **Multiattribute utility theory** deals with utilities that depend on several distinct attributes of states. **Stochastic dominance** is a particularly useful technique for making unambiguous decisions, even without precise utility values for attributes.
- **Decision networks** provide a simple formalism for expressing and solving decision problems. They are a natural extension of Bayesian networks, containing decision and utility nodes in addition to chance nodes.
- Sometimes, solving a problem involves finding more information before making a decision. The **value of information** is defined as the expected improvement in utility compared with making a decision without the information.
- **Expert systems** that incorporate utility information have additional capabilities compared with pure inference systems. In addition to being able to make decisions, they can use the value of information to decide which questions to ask, if any; they can recommend contingency plans; and they can calculate the sensitivity of their decisions to small changes in probability and utility assessments.

---

## BIBLIOGRAPHICAL AND HISTORICAL NOTES

The book *L'art de Penser*, also known as the *Port-Royal Logic* (Arnauld, 1662) states:

To judge what one must do to obtain a good or avoid an evil, it is necessary to consider not only the good and the evil in itself, but also the probability that it happens or does not happen; and to view geometrically the proportion that all these things have together.



Modern texts talk of *utility* rather than good and evil, but this statement correctly notes that one should multiply utility by probability (“view geometrically”) to give expected utility, and maximize that over all outcomes (“all these things”) to “judge what one must do.” It is remarkable how much this got right, 350 years ago, and only 8 years after Pascal and Fermat showed how to use probability correctly. The Port-Royal Logic also marked the first publication of Pascal’s wager.

Daniel Bernoulli (1738), investigating the St. Petersburg paradox (see Exercise 16.3), was the first to realize the importance of preference measurement for lotteries, writing “the *value* of an item must not be based on its *price*, but rather on the *utility* that it yields” (*italics his*). Utilitarian philosopher Jeremy Bentham (1823) proposed the **hedonic calculus** for weighing “pleasures” and “pains,” arguing that all decisions (not just monetary ones) could be reduced to utility comparisons.

The derivation of numerical utilities from preferences was first carried out by Ramsey (1931); the axioms for preference in the present text are closer in form to those rediscovered in *Theory of Games and Economic Behavior* (von Neumann and Morgenstern, 1944). A good presentation of these axioms, in the course of a discussion on risk preference, is given by Howard (1977). Ramsey had derived subjective probabilities (not just utilities) from an agent’s preferences; Savage (1954) and Jeffrey (1983) carry out more recent constructions of this kind. Von Winterfeldt and Edwards (1986) provide a modern perspective on decision analysis and its relationship to human preference structures. The micromort utility measure is discussed by Howard (1989). A 1994 survey by the *Economist* set the value of a life at between \$750,000 and \$2.6 million. However, Richard Thaler (1992) found irrational framing effects on the price one is willing to pay to avoid a risk of death versus the price one is willing to be paid to accept a risk. For a 1/1000 chance, a respondent wouldn’t pay more than \$200 to remove the risk, but wouldn’t accept \$50,000 to take on the risk. How much are people willing to pay for a QALY? When it comes down to a specific case of saving oneself or a family member, the number is approximately “whatever I’ve got.” But we can ask at a societal level: suppose there is a vaccine that would yield  $X$  QALYs but costs  $Y$  dollars; is it worth it? In this case people report a wide range of values from around \$10,000 to \$150,000 per QALY (Prades *et al.*, 2008). QALYs are much more widely used in medical and social policy decision making than are micromorts; see (Russell, 1990) for a typical example of an argument for a major change in public health policy on grounds of increased expected utility measured in QALYs.

The **optimizer’s curse** was brought to the attention of decision analysts in a forceful way by Smith and Winkler (2006), who pointed out that the financial benefits to the client projected by analysts for their proposed course of action almost never materialized. They trace this directly to the bias introduced by selecting an optimal action and show that a more complete Bayesian analysis eliminates the problem. The same underlying concept has been called **post-decision disappointment** by Harrison and March (1984) and was noted in the context of analyzing capital investment projects by Brown (1974). The optimizer’s curse is also closely related to the **winner’s curse** (Capen *et al.*, 1971; Thaler, 1992), which applies to competitive bidding in auctions: whoever wins the auction is very likely to have overestimated the value of the object in question. Capen *et al.* quote a petroleum engineer on the

POST-DECISION  
DISAPPOINTMENT

WINNER’S CURSE

topic of bidding for oil-drilling rights: “If one wins a tract against two or three others he may feel fine about his good fortune. But how should he feel if he won against 50 others? Ill.” Finally, behind both curses is the general phenomenon of **regression to the mean**, whereby individuals selected on the basis of exceptional characteristics previously exhibited will, with high probability, become less exceptional in future.

The Allais paradox, due to Nobel Prize-winning economist Maurice Allais (1953) was tested experimentally (Tversky and Kahneman, 1982; Conlisk, 1989) to show that people are consistently inconsistent in their judgments. The Ellsberg paradox on ambiguity aversion was introduced in the Ph.D. thesis of Daniel Ellsberg (Ellsberg, 1962), who went on to become a military analyst at the RAND Corporation and to leak documents known as The Pentagon Papers, which contributed to the end of the Vietnam war and the resignation of President Nixon. Fox and Tversky (1995) describe a further study of ambiguity aversion. Mark Machina (2005) gives an overview of choice under uncertainty and how it can vary from expected utility theory.

There has been a recent outpouring of more-or-less popular books on human irrationality. The best known is *Predictably Irrational* (Ariely, 2009); others include *Sway* (Brafman and Brafman, 2009), *Nudge* (Thaler and Sunstein, 2009), *Kluge* (Marcus, 2009), *How We Decide* (Lehrer, 2009) and *On Being Certain* (Burton, 2009). They complement the classic (Kahneman *et al.*, 1982) and the article that started it all (Kahneman and Tversky, 1979). The field of evolutionary psychology (Buss, 2005), on the other hand, has run counter to this literature, arguing that humans are quite rational in evolutionarily appropriate contexts. Its adherents point out that irrationality is penalized by definition in an evolutionary context and show that in some cases it is an artifact of the experimental setup (Cummins and Allen, 1998). There has been a recent resurgence of interest in Bayesian models of cognition, overturning decades of pessimism (Oaksford and Chater, 1998; Elio, 2002; Chater and Oaksford, 2008).

Keeney and Raiffa (1976) give a thorough introduction to multiattribute utility theory. They describe early computer implementations of methods for eliciting the necessary parameters for a multiattribute utility function and include extensive accounts of real applications of the theory. In AI, the principal reference for MAUT is Wellman’s (1985) paper, which includes a system called URP (Utility Reasoning Package) that can use a collection of statements about preference independence and conditional independence to analyze the structure of decision problems. The use of stochastic dominance together with qualitative probability models was investigated extensively by Wellman (1988, 1990a). Wellman and Doyle (1992) provide a preliminary sketch of how a complex set of utility-independence relationships might be used to provide a structured model of a utility function, in much the same way that Bayesian networks provide a structured model of joint probability distributions. Bacchus and Grove (1995, 1996) and La Mura and Shoham (1999) give further results along these lines.

Decision theory has been a standard tool in economics, finance, and management science since the 1950s. Until the 1980s, decision trees were the main tool used for representing simple decision problems. Smith (1988) gives an overview of the methodology of decision analysis. Influence diagrams were introduced by Howard and Matheson (1984), based on earlier work at SRI (Miller *et al.*, 1976). Howard and Matheson’s method involved the

derivation of a decision tree from a decision network, but in general the tree is of exponential size. Shachter (1986) developed a method for making decisions based directly on a decision network, without the creation of an intermediate decision tree. This algorithm was also one of the first to provide complete inference for multiply connected Bayesian networks. Zhang *et al.* (1994) showed how to take advantage of conditional independence of information to reduce the size of trees in practice; they use the term *decision network* for networks that use this approach (although others use it as a synonym for influence diagram). Nilsson and Lauritzen (2000) link algorithms for decision networks to ongoing developments in clustering algorithms for Bayesian networks. Koller and Milch (2003) show how influence diagrams can be used to solve games that involve gathering information by opposing players, and Detwarasiti and Shachter (2005) show how influence diagrams can be used as an aid to decision making for a team that shares goals but is unable to share all information perfectly. The collection by Oliver and Smith (1990) has a number of useful articles on decision networks, as does the 1990 special issue of the journal *Networks*. Papers on decision networks and utility modeling also appear regularly in the journals *Management Science* and *Decision Analysis*.

The theory of information value was explored first in the context of statistical experiments, where a quasi-utility (entropy reduction) was used (Lindley, 1956). The Russian control theorist Ruslan Stratonovich (1965) developed the more general theory presented here, in which information has value by virtue of its ability to affect decisions. Stratonovich's work was not known in the West, where Ron Howard (1966) pioneered the same idea. His paper ends with the remark "If information value theory and associated decision theoretic structures do not in the future occupy a large part of the education of engineers, then the engineering profession will find that its traditional role of managing scientific and economic resources for the benefit of man has been forfeited to another profession." To date, the implied revolution in managerial methods has not occurred.

Recent work by Krause and Guestrin (2009) shows that computing the exact non-myopic value of information is intractable even in polytree networks. There are other cases—more restricted than general value of information—in which the myopic algorithm does provide a provably good approximation to the optimal sequence of observations (Krause *et al.*, 2008). In some cases—for example, looking for treasure buried in one of  $n$  places—ranking experiments in order of success probability divided by cost gives an optimal solution (Kadane and Simon, 1977).

Surprisingly few early AI researchers adopted decision-theoretic tools after the early applications in medical decision making described in Chapter 13. One of the few exceptions was Jerry Feldman, who applied decision theory to problems in vision (Feldman and Yakimovsky, 1974) and planning (Feldman and Sproull, 1977). After the resurgence of interest in probabilistic methods in AI in the 1980s, decision-theoretic expert systems gained widespread acceptance (Horvitz *et al.*, 1988; Cowell *et al.*, 2002). In fact, from 1991 onward, the cover design of the journal *Artificial Intelligence* has depicted a decision network, although some artistic license appears to have been taken with the direction of the arrows.

**EXERCISES**

**16.1** (Adapted from David Heckerman.) This exercise concerns the **Almanac Game**, which is used by decision analysts to calibrate numeric estimation. For each of the questions that follow, give your best guess of the answer, that is, a number that you think is as likely to be too high as it is to be too low. Also give your guess at a 25th percentile estimate, that is, a number that you think has a 25% chance of being too high, and a 75% chance of being too low. Do the same for the 75th percentile. (Thus, you should give three estimates in all—low, median, and high—for each question.)

- a. Number of passengers who flew between New York and Los Angeles in 1989.
- b. Population of Warsaw in 1992.
- c. Year in which Coronado discovered the Mississippi River.
- d. Number of votes received by Jimmy Carter in the 1976 presidential election.
- e. Age of the oldest living tree, as of 2002.
- f. Height of the Hoover Dam in feet.
- g. Number of eggs produced in Oregon in 1985.
- h. Number of Buddhists in the world in 1992.
- i. Number of deaths due to AIDS in the United States in 1981.
- j. Number of U.S. patents granted in 1901.

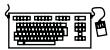
The correct answers appear after the last exercise of this chapter. From the point of view of decision analysis, the interesting thing is not how close your median guesses came to the real answers, but rather how often the real answer came within your 25% and 75% bounds. If it was about half the time, then your bounds are accurate. But if you're like most people, you will be more sure of yourself than you should be, and fewer than half the answers will fall within the bounds. With practice, you can calibrate yourself to give realistic bounds, and thus be more useful in supplying information for decision making. Try this second set of questions and see if there is any improvement:

- a. Year of birth of Zsa Zsa Gabor.
- b. Maximum distance from Mars to the sun in miles.
- c. Value in dollars of exports of wheat from the United States in 1992.
- d. Tons handled by the port of Honolulu in 1991.
- e. Annual salary in dollars of the governor of California in 1993.
- f. Population of San Diego in 1990.
- g. Year in which Roger Williams founded Providence, Rhode Island.
- h. Height of Mt. Kilimanjaro in feet.
- i. Length of the Brooklyn Bridge in feet.
- j. Number of deaths due to automobile accidents in the United States in 1992.

**16.2** Chris considers four used cars before buying the one with maximum expected utility. Pat considers ten cars and does the same. All other things being equal, which one is more likely to have the better car? Which is more likely to be disappointed with their car's quality? By how much (in terms of standard deviations of expected quality)?

**16.3** In 1713, Nicolas Bernoulli stated a puzzle, now called the St. Petersburg paradox, which works as follows. You have the opportunity to play a game in which a fair coin is tossed repeatedly until it comes up heads. If the first heads appears on the  $n$ th toss, you win  $2^n$  dollars.

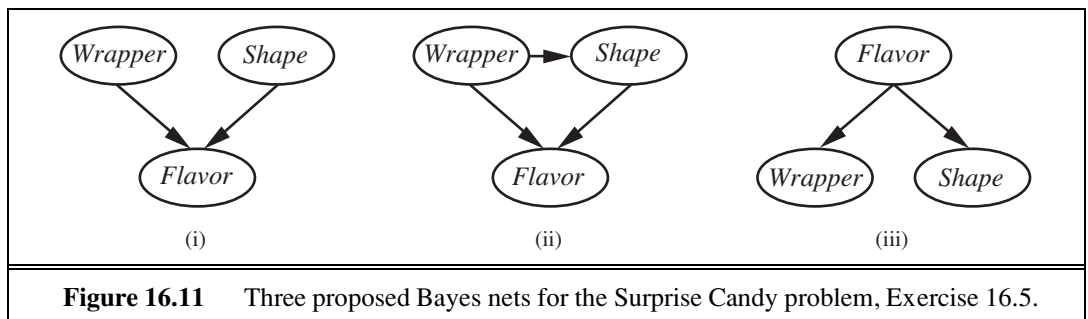
- Show that the expected monetary value of this game is infinite.
- How much would you, personally, pay to play the game?
- Nicolas's cousin Daniel Bernoulli resolved the apparent paradox in 1738 by suggesting that the utility of money is measured on a logarithmic scale (i.e.,  $U(S_n) = a \log_2 n + b$ , where  $S_n$  is the state of having  $\$n$ ). What is the expected utility of the game under this assumption?
- What is the maximum amount that it would be rational to pay to play the game, assuming that one's initial wealth is  $\$k$ ?



**16.4** Write a computer program to automate the process in Exercise 16.9. Try your program out on several people of different net worth and political outlook. Comment on the consistency of your results, both for an individual and across individuals.

**16.5** The Surprise Candy Company makes candy in two flavors: 70% are strawberry flavor and 30% are anchovy flavor. Each new piece of candy starts out with a round shape; as it moves along the production line, a machine randomly selects a certain percentage to be trimmed into a square; then, each piece is wrapped in a wrapper whose color is chosen randomly to be red or brown. 80% of the strawberry candies are round and 80% have a red wrapper, while 90% of the anchovy candies are square and 90% have a brown wrapper. All candies are sold individually in sealed, identical, black boxes.

Now you, the customer, have just bought a Surprise candy at the store but have not yet opened the box. Consider the three Bayes nets in Figure 16.11.



**Figure 16.11** Three proposed Bayes nets for the Surprise Candy problem, Exercise 16.5.

- Which network(s) can correctly represent  $\mathbf{P}(\text{Flavor}, \text{Wrapper}, \text{Shape})$ ?
- Which network is the best representation for this problem?

- c. Does network (i) assert that  $\mathbf{P}(\text{Wrapper}|\text{Shape}) = \mathbf{P}(\text{Wrapper})$ ?
- d. What is the probability that your candy has a red wrapper?
- e. In the box is a round candy with a red wrapper. What is the probability that its flavor is strawberry?
- f. A unwrapped strawberry candy is worth  $s$  on the open market and an unwrapped anchovy candy is worth  $a$ . Write an expression for the value of an unopened candy box.
- g. A new law prohibits trading of unwrapped candies, but it is still legal to trade wrapped candies (out of the box). Is an unopened candy box now worth more than less than, or the same as before?

**16.6** Prove that the judgments  $B \succ A$  and  $C \succ D$  in the Allais paradox (page 620) violate the axiom of substitutability.

**16.7** Consider the Allais paradox described on page 620: an agent who prefers  $B$  over  $A$  (taking the sure thing), and  $C$  over  $D$  (taking the higher EMV) is not acting rationally, according to utility theory. Do you think this indicates a problem for the agent, a problem for the theory, or no problem at all? Explain.

**16.8** Tickets to a lottery cost \$1. There are two possible prizes: a \$10 payoff with probability  $1/50$ , and a \$1,000,000 payoff with probability  $1/2,000,000$ . What is the expected monetary value of a lottery ticket? When (if ever) is it rational to buy a ticket? Be precise—show an equation involving utilities. You may assume current wealth of  $\$k$  and that  $U(S_k) = 0$ . You may also assume that  $U(S_{k+10}) = 10 \times U(S_{k+1})$ , but you may not make any assumptions about  $U(S_{k+1,000,000})$ . Sociological studies show that people with lower income buy a disproportionate number of lottery tickets. Do you think this is because they are worse decision makers or because they have a different utility function? Consider the value of contemplating the possibility of winning the lottery versus the value of contemplating becoming an action hero while watching an adventure movie.

**16.9** Assess your own utility for different incremental amounts of money by running a series of preference tests between some definite amount  $M_1$  and a lottery  $[p, M_2; (1-p), 0]$ . Choose different values of  $M_1$  and  $M_2$ , and vary  $p$  until you are indifferent between the two choices. Plot the resulting utility function.

**16.10** How much is a micromort worth to you? Devise a protocol to determine this. Ask questions based both on paying to avoid risk and being paid to accept risk.

**16.11** Let continuous variables  $X_1, \dots, X_k$  be independently distributed according to the same probability density function  $f(x)$ . Prove that the density function for  $\max\{X_1, \dots, X_k\}$  is given by  $kf(x)(F(x))^{k-1}$ , where  $F$  is the cumulative distribution for  $f$ .

**16.12** Economists often make use of an exponential utility function for money:  $U(x) = -e^{x/R}$ , where  $R$  is a positive constant representing an individual's risk tolerance. Risk tolerance reflects how likely an individual is to accept a lottery with a particular expected monetary value (EMV) versus some certain payoff. As  $R$  (which is measured in the same units as  $x$ ) becomes larger, the individual becomes less risk-averse.

- a. Assume Mary has an exponential utility function with  $R = \$500$ . Mary is given the choice between receiving \$500 with certainty (probability 1) or participating in a lottery which has a 60% probability of winning \$5000 and a 40% probability of winning nothing. Assuming Mary acts rationally, which option would she choose? Show how you derived your answer.
- b. Consider the choice between receiving \$100 with certainty (probability 1) or participating in a lottery which has a 50% probability of winning \$500 and a 50% probability of winning nothing. Approximate the value of  $R$  (to 3 significant digits) in an exponential utility function that would cause an individual to be indifferent to these two alternatives. (You might find it helpful to write a short program to help you solve this problem.)

**16.13** Repeat Exercise 16.16, using the action-utility representation shown in Figure 16.7.

**16.14** For either of the airport-siting diagrams from Exercises 16.16 and 16.13, to which conditional probability table entry is the utility most sensitive, given the available evidence?

**16.15** Consider a student who has the choice to buy or not buy a textbook for a course. We'll model this as a decision problem with one Boolean decision node,  $B$ , indicating whether the agent chooses to buy the book, and two Boolean chance nodes,  $M$ , indicating whether the student has mastered the material in the book, and  $P$ , indicating whether the student passes the course. Of course, there is also a utility node,  $U$ . A certain student, Sam, has an additive utility function: 0 for not buying the book and -\$100 for buying it; and \$2000 for passing the course and 0 for not passing. Sam's conditional probability estimates are as follows:

$$\begin{aligned} P(p|b, m) &= 0.9 & P(m|b) &= 0.9 \\ P(p|b, \neg m) &= 0.5 & P(m|\neg b) &= 0.7 \\ P(p|\neg b, m) &= 0.8 \\ P(p|\neg b, \neg m) &= 0.3 \end{aligned}$$

You might think that  $P$  would be independent of  $B$  given  $M$ . But this course has an open-book final—so having the book helps.

- a. Draw the decision network for this problem.
- b. Compute the expected utility of buying the book and of not buying it.
- c. What should Sam do?



**16.16** This exercise completes the analysis of the airport-siting problem in Figure 16.6.

- a. Provide reasonable variable domains, probabilities, and utilities for the network, assuming that there are three possible sites.
- b. Solve the decision problem.
- c. What happens if changes in technology mean that each aircraft generates half the noise?
- d. What if noise avoidance becomes three times more important?
- e. Calculate the VPI for *AirTraffic*, *Litigation*, and *Construction* in your model.

**16.17** (Adapted from Pearl (1988).) A used-car buyer can decide to carry out various tests with various costs (e.g., kick the tires, take the car to a qualified mechanic) and then, depending on the outcome of the tests, decide which car to buy. We will assume that the buyer is deciding whether to buy car  $c_1$ , that there is time to carry out at most one test, and that  $t_1$  is the test of  $c_1$  and costs \$50.

A car can be in good shape (quality  $q^+$ ) or bad shape (quality  $q^-$ ), and the tests might help indicate what shape the car is in. Car  $c_1$  costs \$1,500, and its market value is \$2,000 if it is in good shape; if not, \$700 in repairs will be needed to make it in good shape. The buyer's estimate is that  $c_1$  has a 70% chance of being in good shape.

- a. Draw the decision network that represents this problem.
- b. Calculate the expected net gain from buying  $c_1$ , given no test.
- c. Tests can be described by the probability that the car will pass or fail the test given that the car is in good or bad shape. We have the following information:  

$$P(\text{pass}(c_1, t_1) | q^+(c_1)) = 0.8$$

$$P(\text{pass}(c_1, t_1) | q^-(c_1)) = 0.35$$
 Use Bayes' theorem to calculate the probability that the car will pass (or fail) its test and hence the probability that it is in good (or bad) shape given each possible test outcome.
- d. Calculate the optimal decisions given either a pass or a fail, and their expected utilities.
- e. Calculate the value of information of the test, and derive an optimal conditional plan for the buyer.

**16.18** Recall the definition of *value of information* in Section 16.6.

- a. Prove that the value of information is nonnegative and order independent.
- b. Explain why it is that some people would prefer not to get some information—for example, not wanting to know the sex of their baby when an ultrasound is done.
- c. A function  $f$  on sets is **submodular** if, for any element  $x$  and any sets  $A$  and  $B$  such that  $A \subseteq B$ , adding  $x$  to  $A$  gives a greater increase in  $f$  than adding  $x$  to  $B$ :

$$A \subseteq B \Rightarrow (f(A \cup \{x\}) - f(A)) \geq (f(B \cup \{x\}) - f(B)) .$$

Submodularity captures the intuitive notion of *diminishing returns*. Is the value of information, viewed as a function  $f$  on sets of possible observations, submodular? Prove this or find a counterexample.

SUBMODULARITY

The answers to Exercise 16.1 (where M stands for million): First set: 3M, 1.6M, 1541, 41M, 4768, 221, 649M, 295M, 132, 25,546. Second set: 1917, 155M, 4,500M, 11M, 120,000, 1.1M, 1636, 19,340, 1,595, 41,710.



# 17 MAKING COMPLEX DECISIONS

*In which we examine methods for deciding what to do today, given that we may decide again tomorrow.*

SEQUENTIAL  
DECISION PROBLEM

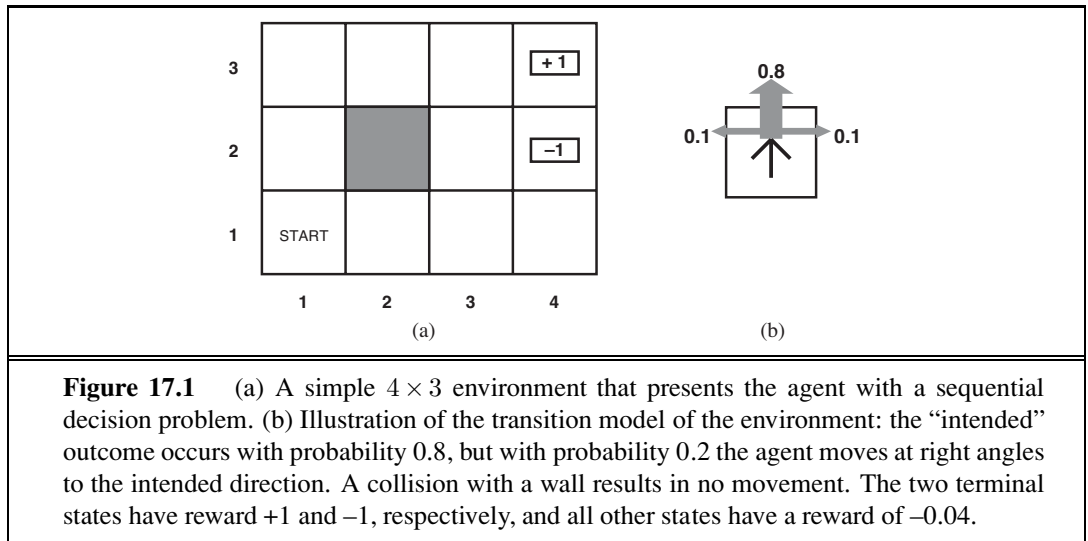
In this chapter, we address the computational issues involved in making decisions in a stochastic environment. Whereas Chapter 16 was concerned with one-shot or episodic decision problems, in which the utility of each action's outcome was well known, we are concerned here with **sequential decision problems**, in which the agent's utility depends on a sequence of decisions. Sequential decision problems incorporate utilities, uncertainty, and sensing, and include search and planning problems as special cases. Section 17.1 explains how sequential decision problems are defined, and Sections 17.2 and 17.3 explain how they can be solved to produce optimal behavior that balances the risks and rewards of acting in an uncertain environment. Section 17.4 extends these ideas to the case of partially observable environments, and Section 17.4.3 develops a complete design for decision-theoretic agents in partially observable environments, combining dynamic Bayesian networks from Chapter 15 with decision networks from Chapter 16.

The second part of the chapter covers environments with multiple agents. In such environments, the notion of optimal behavior is complicated by the interactions among the agents. Section 17.5 introduces the main ideas of **game theory**, including the idea that rational agents might need to behave randomly. Section 17.6 looks at how multiagent systems can be designed so that multiple agents can achieve a common goal.

## 17.1 SEQUENTIAL DECISION PROBLEMS

---

Suppose that an agent is situated in the  $4 \times 3$  environment shown in Figure 17.1(a). Beginning in the start state, it must choose an action at each time step. The interaction with the environment terminates when the agent reaches one of the goal states, marked +1 or -1. Just as for search problems, the actions available to the agent in each state are given by  $\text{ACTIONS}(s)$ , sometimes abbreviated to  $A(s)$ ; in the  $4 \times 3$  environment, the actions in every state are *Up*, *Down*, *Left*, and *Right*. We assume for now that the environment is **fully observable**, so that the agent always knows where it is.



If the environment were deterministic, a solution would be easy: *[Up, Up, Right, Right, Right]*. Unfortunately, the environment won’t always go along with this solution, because the actions are unreliable. The particular model of stochastic motion that we adopt is illustrated in Figure 17.1(b). Each action achieves the intended effect with probability 0.8, but the rest of the time, the action moves the agent at right angles to the intended direction. Furthermore, if the agent bumps into a wall, it stays in the same square. For example, from the start square (1,1), the action *Up* moves the agent to (1,2) with probability 0.8, but with probability 0.1, it moves right to (2,1), and with probability 0.1, it moves left, bumps into the wall, and stays in (1,1). In such an environment, the sequence *[Up, Up, Right, Right, Right]* goes up around the barrier and reaches the goal state at (4,3) with probability  $0.8^5 = 0.32768$ . There is also a small chance of accidentally reaching the goal by going the other way around with probability  $0.1^4 \times 0.8$ , for a grand total of 0.32776. (See also Exercise 17.1.)

As in Chapter 3, the **transition model** (or just “model,” whenever no confusion can arise) describes the outcome of each action in each state. Here, the outcome is stochastic, so we write  $P(s' | s, a)$  to denote the probability of reaching state  $s'$  if action  $a$  is done in state  $s$ . We will assume that transitions are **Markovian** in the sense of Chapter 15, that is, the probability of reaching  $s'$  from  $s$  depends only on  $s$  and not on the history of earlier states. For now, you can think of  $P(s' | s, a)$  as a big three-dimensional table containing probabilities. Later, in Section 17.4.3, we will see that the transition model can be represented as a **dynamic Bayesian network**, just as in Chapter 15.

To complete the definition of the task environment, we must specify the utility function for the agent. Because the decision problem is sequential, the utility function will depend on a sequence of states—an **environment history**—rather than on a single state. Later in this section, we investigate how such utility functions can be specified in general; for now, we simply stipulate that in each state  $s$ , the agent receives a **reward**  $R(s)$ , which may be positive or negative, but must be bounded. For our particular example, the reward is -0.04 in all states except the terminal states (which have rewards +1 and -1). The utility of an

environment history is just (for now) the *sum* of the rewards received. For example, if the agent reaches the +1 state after 10 steps, its total utility will be 0.6. The negative reward of  $-0.04$  gives the agent an incentive to reach (4,3) quickly, so our environment is a stochastic generalization of the search problems of Chapter 3. Another way of saying this is that the agent does not enjoy living in this environment and so wants to leave as soon as possible.

MARKOV DECISION  
PROCESS

To sum up: a sequential decision problem for a fully observable, stochastic environment with a Markovian transition model and additive rewards is called a **Markov decision process**, or **MDP**, and consists of a set of states (with an initial state  $s_0$ ); a set  $\text{ACTIONS}(s)$  of actions in each state; a transition model  $P(s' | s, a)$ ; and a reward function  $R(s)$ .<sup>1</sup>

POLICY

The next question is, what does a solution to the problem look like? We have seen that any fixed action sequence won't solve the problem, because the agent might end up in a state other than the goal. Therefore, a solution must specify what the agent should do for *any* state that the agent might reach. A solution of this kind is called a **policy**. It is traditional to denote a policy by  $\pi$ , and  $\pi(s)$  is the action recommended by the policy  $\pi$  for state  $s$ . If the agent has a complete policy, then no matter what the outcome of any action, the agent will always know what to do next.

OPTIMAL POLICY

Each time a given policy is executed starting from the initial state, the stochastic nature of the environment may lead to a different environment history. The quality of a policy is therefore measured by the *expected* utility of the possible environment histories generated by that policy. An **optimal policy** is a policy that yields the highest expected utility. We use  $\pi^*$  to denote an optimal policy. Given  $\pi^*$ , the agent decides what to do by consulting its current percept, which tells it the current state  $s$ , and then executing the action  $\pi^*(s)$ . A policy represents the agent function explicitly and is therefore a description of a simple reflex agent, computed from the information used for a utility-based agent.

An optimal policy for the world of Figure 17.1 is shown in Figure 17.2(a). Notice that, because the cost of taking a step is fairly small compared with the penalty for ending up in (4,2) by accident, the optimal policy for the state (3,1) is conservative. The policy recommends taking the long way round, rather than taking the shortcut and thereby risking entering (4,2).

The balance of risk and reward changes depending on the value of  $R(s)$  for the nonterminal states. Figure 17.2(b) shows optimal policies for four different ranges of  $R(s)$ . When  $R(s) \leq -1.6284$ , life is so painful that the agent heads straight for the nearest exit, even if the exit is worth  $-1$ . When  $-0.4278 \leq R(s) \leq -0.0850$ , life is quite unpleasant; the agent takes the shortest route to the +1 state and is willing to risk falling into the  $-1$  state by accident. In particular, the agent takes the shortcut from (3,1). When life is only slightly dreary ( $-0.0221 < R(s) < 0$ ), the optimal policy takes *no risks at all*. In (4,1) and (3,2), the agent heads directly away from the  $-1$  state so that it cannot fall in by accident, even though this means banging its head against the wall quite a few times. Finally, if  $R(s) > 0$ , then life is positively enjoyable and the agent avoids *both* exits. As long as the actions in (4,1), (3,2),

<sup>1</sup> Some definitions of MDPs allow the reward to depend on the action and outcome too, so the reward function is  $R(s, a, s')$ . This simplifies the description of some environments but does not change the problem in any fundamental way, as shown in Exercise 17.4.



NONSTATIONARY  
POLICY

STATIONARY POLICY

*the optimal action in a given state could change over time.* We say that the optimal policy for a finite horizon is **nonstationary**. With no fixed time limit, on the other hand, there is no reason to behave differently in the same state at different times. Hence, the optimal action depends only on the current state, and the optimal policy is **stationary**. Policies for the infinite-horizon case are therefore simpler than those for the finite-horizon case, and we deal mainly with the infinite-horizon case in this chapter. (We will see later that for partially observable environments, the infinite-horizon case is not so simple.) Note that “infinite horizon” does not necessarily mean that all state sequences are infinite; it just means that there is no fixed deadline. In particular, there can be finite state sequences in an infinite-horizon MDP containing a terminal state.

STATIONARY  
PREFERENCE

The next question we must decide is how to calculate the utility of state sequences. In the terminology of multiattribute utility theory, each state  $s_i$  can be viewed as an **attribute** of the state sequence  $[s_0, s_1, s_2, \dots]$ . To obtain a simple expression in terms of the attributes, we will need to make some sort of preference-independence assumption. The most natural assumption is that the agent’s preferences between state sequences are **stationary**. Stationarity for preferences means the following: if two state sequences  $[s_0, s_1, s_2, \dots]$  and  $[s'_0, s'_1, s'_2, \dots]$  begin with the same state (i.e.,  $s_0 = s'_0$ ), then the two sequences should be preference-ordered the same way as the sequences  $[s_1, s_2, \dots]$  and  $[s'_1, s'_2, \dots]$ . In English, this means that if you prefer one future to another starting tomorrow, then you should still prefer that future if it were to start today instead. Stationarity is a fairly innocuous-looking assumption with very strong consequences: it turns out that under stationarity there are just two coherent ways to assign utilities to sequences:

ADDITIVE REWARD

1. **Additive rewards:** The utility of a state sequence is

$$U_h([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$$

The  $4 \times 3$  world in Figure 17.1 uses additive rewards. Notice that additivity was used implicitly in our use of path cost functions in heuristic search algorithms (Chapter 3).

DISCOUNTED  
REWARD

2. **Discounted rewards:** The utility of a state sequence is

$$U_h([s_0, s_1, s_2, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots,$$

DISCOUNT FACTOR

where the **discount factor**  $\gamma$  is a number between 0 and 1. The discount factor describes the preference of an agent for current rewards over future rewards. When  $\gamma$  is close to 0, rewards in the distant future are viewed as insignificant. When  $\gamma$  is 1, discounted rewards are exactly equivalent to additive rewards, so additive rewards are a special case of discounted rewards. Discounting appears to be a good model of both animal and human preferences over time. A discount factor of  $\gamma$  is equivalent to an interest rate of  $(1/\gamma) - 1$ .

For reasons that will shortly become clear, we assume discounted rewards in the remainder of the chapter, although sometimes we allow  $\gamma = 1$ .

Lurking beneath our choice of infinite horizons is a problem: if the environment does not contain a terminal state, or if the agent never reaches one, then all environment histories will be infinitely long, and utilities with additive, undiscounted rewards will generally be

infinite. While we can agree that  $+\infty$  is better than  $-\infty$ , comparing two state sequences with  $+\infty$  utility is more difficult. There are three solutions, two of which we have seen already:

1. With discounted rewards, the utility of an infinite sequence is *finite*. In fact, if  $\gamma < 1$  and rewards are bounded by  $\pm R_{\max}$ , we have

$$U_h([s_0, s_1, s_2, \dots]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq \sum_{t=0}^{\infty} \gamma^t R_{\max} = R_{\max}/(1 - \gamma), \quad (17.1)$$

using the standard formula for the sum of an infinite geometric series.

2. If the environment contains terminal states *and if the agent is guaranteed to get to one eventually*, then we will never need to compare infinite sequences. A policy that is guaranteed to reach a terminal state is called a **proper policy**. With proper policies, we can use  $\gamma = 1$  (i.e., additive rewards). The first three policies shown in Figure 17.2(b) are proper, but the fourth is improper. It gains infinite total reward by staying away from the terminal states when the reward for the nonterminal states is positive. The existence of improper policies can cause the standard algorithms for solving MDPs to fail with additive rewards, and so provides a good reason for using discounted rewards.

PROPER POLICY

3. Infinite sequences can be compared in terms of the **average reward** obtained per time step. Suppose that square (1,1) in the  $4 \times 3$  world has a reward of 0.1 while the other nonterminal states have a reward of 0.01. Then a policy that does its best to stay in (1,1) will have higher average reward than one that stays elsewhere. Average reward is a useful criterion for some problems, but the analysis of average-reward algorithms is beyond the scope of this book.

AVERAGE REWARD

In sum, discounted rewards present the fewest difficulties in evaluating state sequences.

### 17.1.2 Optimal policies and the utilities of states

Having decided that the utility of a given state sequence is the sum of discounted rewards obtained during the sequence, we can compare policies by comparing the *expected* utilities obtained when executing them. We assume the agent is in some initial state  $s$  and define  $S_t$  (a random variable) to be the state the agent reaches at time  $t$  when executing a particular policy  $\pi$ . (Obviously,  $S_0 = s$ , the state the agent is in now.) The probability distribution over state sequences  $S_1, S_2, \dots$ , is determined by the initial state  $s$ , the policy  $\pi$ , and the transition model for the environment.

The expected utility obtained by executing  $\pi$  starting in  $s$  is given by

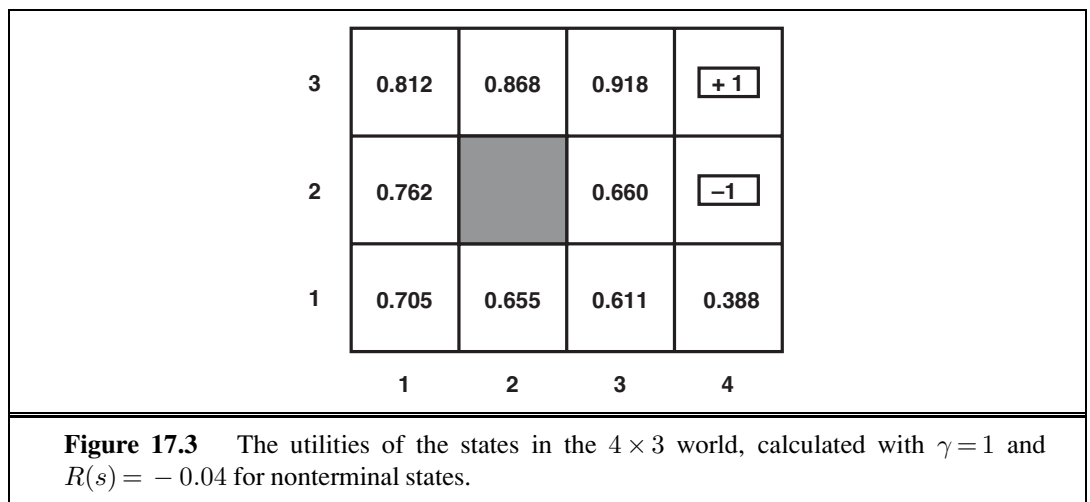
$$U^\pi(s) = E \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t) \right], \quad (17.2)$$

where the expectation is with respect to the probability distribution over state sequences determined by  $s$  and  $\pi$ . Now, out of all the policies the agent could choose to execute starting in  $s$ , one (or more) will have higher expected utilities than all the others. We'll use  $\pi_s^*$  to denote one of these policies:

$$\pi_s^* = \operatorname{argmax}_{\pi} U^\pi(s). \quad (17.3)$$

Remember that  $\pi_s^*$  is a policy, so it recommends an action for every state; its connection with  $s$  in particular is that it's an optimal policy when  $s$  is the starting state. A remarkable consequence of using discounted utilities with infinite horizons is that the optimal policy is *independent* of the starting state. (Of course, the *action sequence* won't be independent; remember that a policy is a function specifying an action for each state.) This fact seems intuitively obvious: if policy  $\pi_a^*$  is optimal starting in  $a$  and policy  $\pi_b^*$  is optimal starting in  $b$ , then, when they reach a third state  $c$ , there's no good reason for them to disagree with each other, or with  $\pi_c^*$ , about what to do next.<sup>2</sup> So we can simply write  $\pi^*$  for an optimal policy.

Given this definition, the true utility of a state is just  $U^{\pi^*}(s)$ —that is, the expected sum of discounted rewards if the agent executes an optimal policy. We write this as  $U(s)$ , matching the notation used in Chapter 16 for the utility of an outcome. Notice that  $U(s)$  and  $R(s)$  are quite different quantities;  $R(s)$  is the “short term” reward for being in  $s$ , whereas  $U(s)$  is the “long term” total reward from  $s$  onward. Figure 17.3 shows the utilities for the  $4 \times 3$  world. Notice that the utilities are higher for states closer to the +1 exit, because fewer steps are required to reach the exit.



The utility function  $U(s)$  allows the agent to select actions by using the principle of maximum expected utility from Chapter 16—that is, choose the action that maximizes the expected utility of the subsequent state:

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s' | s, a) U(s'). \quad (17.4)$$

The next two sections describe algorithms for finding optimal policies.

<sup>2</sup> Although this seems obvious, it does not hold for finite-horizon policies or for other ways of combining rewards over time. The proof follows directly from the uniqueness of the utility function on states, as shown in Section 17.2.

## 17.2 VALUE ITERATION

### VALUE ITERATION

In this section, we present an algorithm, called **value iteration**, for calculating an optimal policy. The basic idea is to calculate the utility of each state and then use the state utilities to select an optimal action in each state.

### 17.2.1 The Bellman equation for utilities



Section 17.1.2 defined the utility of being in a state as the expected sum of discounted rewards from that point onwards. From this, it follows that there is a direct relationship between the utility of a state and the utility of its neighbors: *the utility of a state is the immediate reward for that state plus the expected discounted utility of the next state, assuming that the agent chooses the optimal action.* That is, the utility of a state is given by

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U(s'). \quad (17.5)$$

### BELLMAN EQUATION

This is called the **Bellman equation**, after Richard Bellman (1957). The utilities of the states—defined by Equation (17.2) as the expected utility of subsequent state sequences—are solutions of the set of Bellman equations. In fact, they are the *unique* solutions, as we show in Section 17.2.3.

Let us look at one of the Bellman equations for the  $4 \times 3$  world. The equation for the state (1,1) is

$$U(1,1) = -0.04 + \gamma \max \begin{cases} 0.8U(1,2) + 0.1U(2,1) + 0.1U(1,1), & (Up) \\ 0.9U(1,1) + 0.1U(1,2), & (Left) \\ 0.9U(1,1) + 0.1U(2,1), & (Down) \\ 0.8U(2,1) + 0.1U(1,2) + 0.1U(1,1) \end{cases} \quad (Right)$$

When we plug in the numbers from Figure 17.3, we find that *Up* is the best action.

### 17.2.2 The value iteration algorithm

The Bellman equation is the basis of the value iteration algorithm for solving MDPs. If there are  $n$  possible states, then there are  $n$  Bellman equations, one for each state. The  $n$  equations contain  $n$  unknowns—the utilities of the states. So we would like to solve these simultaneous equations to find the utilities. There is one problem: the equations are *nonlinear*, because the “max” operator is not a linear operator. Whereas systems of linear equations can be solved quickly using linear algebra techniques, systems of nonlinear equations are more problematic. One thing to try is an *iterative* approach. We start with arbitrary initial values for the utilities, calculate the right-hand side of the equation, and plug it into the left-hand side—thereby updating the utility of each state from the utilities of its neighbors. We repeat this until we reach an equilibrium. Let  $U_i(s)$  be the utility value for state  $s$  at the  $i$ th iteration. The iteration step, called a **Bellman update**, looks like this:

### BELLMAN UPDATE

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U_i(s'), \quad (17.6)$$



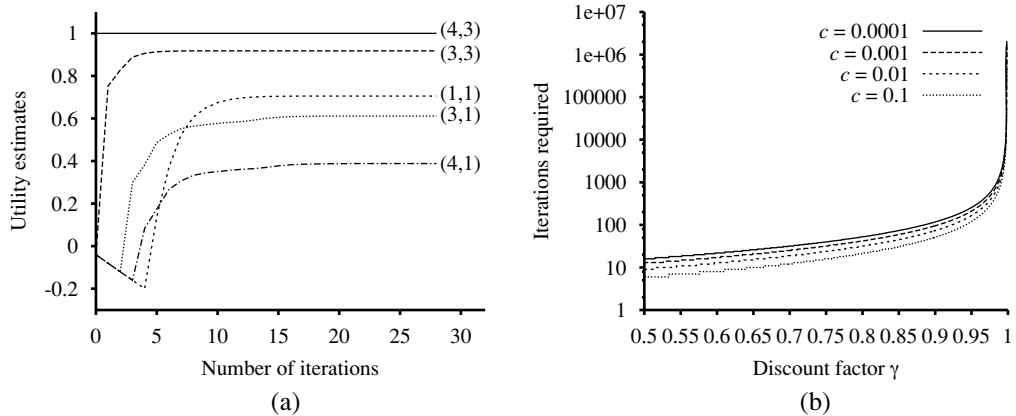
```

function VALUE-ITERATION( $mdp, \epsilon$ ) returns a utility function
  inputs:  $mdp$ , an MDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ ,
           rewards  $R(s)$ , discount  $\gamma$ 
            $\epsilon$ , the maximum error allowed in the utility of any state
  local variables:  $U, U'$ , vectors of utilities for states in  $S$ , initially zero
                     $\delta$ , the maximum change in the utility of any state in an iteration

  repeat
     $U \leftarrow U'; \delta \leftarrow 0$ 
    for each state  $s$  in  $S$  do
       $U'[s] \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s']$ 
      if  $|U'[s] - U[s]| > \delta$  then  $\delta \leftarrow |U'[s] - U[s]|$ 
  until  $\delta < \epsilon(1 - \gamma)/\gamma$ 
  return  $U$ 

```

**Figure 17.4** The value iteration algorithm for calculating utilities of states. The termination condition is from Equation (17.8).



**Figure 17.5** (a) Graph showing the evolution of the utilities of selected states using value iteration. (b) The number of value iterations  $k$  required to guarantee an error of at most  $\epsilon = c \cdot R_{\max}$ , for different values of  $c$ , as a function of the discount factor  $\gamma$ .

where the update is assumed to be applied simultaneously to all the states at each iteration. If we apply the Bellman update infinitely often, we are guaranteed to reach an equilibrium (see Section 17.2.3), in which case the final utility values must be solutions to the Bellman equations. In fact, they are also the *unique* solutions, and the corresponding policy (obtained using Equation (17.4)) is optimal. The algorithm, called VALUE-ITERATION, is shown in Figure 17.4.

We can apply value iteration to the  $4 \times 3$  world in Figure 17.1(a). Starting with initial values of zero, the utilities evolve as shown in Figure 17.5(a). Notice how the states at differ-

ent distances from (4,3) accumulate negative reward until a path is found to (4,3), whereupon the utilities start to increase. We can think of the value iteration algorithm as *propagating information* through the state space by means of local updates.

### 17.2.3 Convergence of value iteration

We said that value iteration eventually converges to a unique set of solutions of the Bellman equations. In this section, we explain why this happens. We introduce some useful mathematical ideas along the way, and we obtain some methods for assessing the error in the utility function returned when the algorithm is terminated early; this is useful because it means that we don't have to run forever. This section is quite technical.

CONTRACTION

The basic concept used in showing that value iteration converges is the notion of a **contraction**. Roughly speaking, a contraction is a function of one argument that, when applied to two different inputs in turn, produces two output values that are “closer together,” by at least some constant factor, than the original inputs. For example, the function “divide by two” is a contraction, because, after we divide any two numbers by two, their difference is halved. Notice that the “divide by two” function has a fixed point, namely zero, that is unchanged by the application of the function. From this example, we can discern two important properties of contractions:

- A contraction has only one fixed point; if there were two fixed points they would not get closer together when the function was applied, so it would not be a contraction.
- When the function is applied to any argument, the value must get closer to the fixed point (because the fixed point does not move), so repeated application of a contraction always reaches the fixed point in the limit.

Now, suppose we view the Bellman update (Equation (17.6)) as an operator  $B$  that is applied simultaneously to update the utility of every state. Let  $U_i$  denote the vector of utilities for all the states at the  $i$ th iteration. Then the Bellman update equation can be written as

$$U_{i+1} \leftarrow B U_i .$$

MAX NORM

Next, we need a way to measure distances between utility vectors. We will use the **max norm**, which measures the “length” of a vector by the absolute value of its biggest component:

$$\|U\| = \max_s |U(s)| .$$

With this definition, the “distance” between two vectors,  $\|U - U'\|$ , is the maximum difference between any two corresponding elements. The main result of this section is the following: *Let  $U_i$  and  $U'_i$  be any two utility vectors. Then we have*

$$\|B U_i - B U'_i\| \leq \gamma \|U_i - U'_i\| . \quad (17.7)$$

*That is, the Bellman update is a contraction by a factor of  $\gamma$  on the space of utility vectors.* (Exercise 17.6 provides some guidance on proving this claim.) Hence, from the properties of contractions in general, it follows that value iteration always converges to a unique solution of the Bellman equations whenever  $\gamma < 1$ .



We can also use the contraction property to analyze the *rate* of convergence to a solution. In particular, we can replace  $U'_i$  in Equation (17.7) with the *true* utilities  $U$ , for which  $BU = U$ . Then we obtain the inequality

$$\|BU_i - U\| \leq \gamma \|U_i - U\|.$$

So, if we view  $\|U_i - U\|$  as the *error* in the estimate  $U_i$ , we see that the error is reduced by a factor of at least  $\gamma$  on each iteration. This means that value iteration converges exponentially fast. We can calculate the number of iterations required to reach a specified error bound  $\epsilon$  as follows: First, recall from Equation (17.1) that the utilities of all states are bounded by  $\pm R_{\max}/(1 - \gamma)$ . This means that the maximum initial error  $\|U_0 - U\| \leq 2R_{\max}/(1 - \gamma)$ . Suppose we run for  $N$  iterations to reach an error of at most  $\epsilon$ . Then, because the error is reduced by at least  $\gamma$  each time, we require  $\gamma^N \cdot 2R_{\max}/(1 - \gamma) \leq \epsilon$ . Taking logs, we find

$$N = \lceil \log(2R_{\max}/\epsilon(1 - \gamma)) / \log(1/\gamma) \rceil$$

iterations suffice. Figure 17.5(b) shows how  $N$  varies with  $\gamma$ , for different values of the ratio  $\epsilon/R_{\max}$ . The good news is that, because of the exponentially fast convergence,  $N$  does not depend much on the ratio  $\epsilon/R_{\max}$ . The bad news is that  $N$  grows rapidly as  $\gamma$  becomes close to 1. We can get fast convergence if we make  $\gamma$  small, but this effectively gives the agent a short horizon and could miss the long-term effects of the agent's actions.

The error bound in the preceding paragraph gives some idea of the factors influencing the run time of the algorithm, but is sometimes overly conservative as a method of deciding when to stop the iteration. For the latter purpose, we can use a bound relating the error to the size of the Bellman update on any given iteration. From the contraction property (Equation (17.7)), it can be shown that if the update is small (i.e., no state's utility changes by much), then the error, compared with the true utility function, also is small. More precisely,

$$\text{if } \|U_{i+1} - U_i\| < \epsilon(1 - \gamma)/\gamma \text{ then } \|U_{i+1} - U\| < \epsilon. \quad (17.8)$$

This is the termination condition used in the VALUE-ITERATION algorithm of Figure 17.4.

So far, we have analyzed the error in the utility function returned by the value iteration algorithm. *What the agent really cares about, however, is how well it will do if it makes its decisions on the basis of this utility function.* Suppose that after  $i$  iterations of value iteration, the agent has an estimate  $U_i$  of the true utility  $U$  and obtains the MEU policy  $\pi_i$  based on one-step look-ahead using  $U_i$  (as in Equation (17.4)). Will the resulting behavior be nearly as good as the optimal behavior? This is a crucial question for any real agent, and it turns out that the answer is yes.  $U^{\pi_i}(s)$  is the utility obtained if  $\pi_i$  is executed starting in  $s$ , and the **policy loss**  $\|U^{\pi_i} - U\|$  is the most the agent can lose by executing  $\pi_i$  instead of the optimal policy  $\pi^*$ . The policy loss of  $\pi_i$  is connected to the error in  $U_i$  by the following inequality:

$$\text{if } \|U_i - U\| < \epsilon \text{ then } \|U^{\pi_i} - U\| < 2\epsilon\gamma/(1 - \gamma). \quad (17.9)$$

In practice, it often occurs that  $\pi_i$  becomes optimal long before  $U_i$  has converged. Figure 17.6 shows how the maximum error in  $U_i$  and the policy loss approach zero as the value iteration process proceeds for the  $4 \times 3$  environment with  $\gamma = 0.9$ . The policy  $\pi_i$  is optimal when  $i = 4$ , even though the maximum error in  $U_i$  is still 0.46.

Now we have everything we need to use value iteration in practice. We know that it converges to the correct utilities, we can bound the error in the utility estimates if we



POLICY LOSS

stop after a finite number of iterations, and we can bound the policy loss that results from executing the corresponding MEU policy. As a final note, all of the results in this section depend on discounting with  $\gamma < 1$ . If  $\gamma = 1$  and the environment contains terminal states, then a similar set of convergence results and error bounds can be derived whenever certain technical conditions are satisfied.

### 17.3 POLICY ITERATION

In the previous section, we observed that it is possible to get an optimal policy even when the utility function estimate is inaccurate. If one action is clearly better than all others, then the exact magnitude of the utilities on the states involved need not be precise. This insight suggests an alternative way to find optimal policies. The **policy iteration** algorithm alternates the following two steps, beginning from some initial policy  $\pi_0$ :

POLICY ITERATION

POLICY EVALUATION

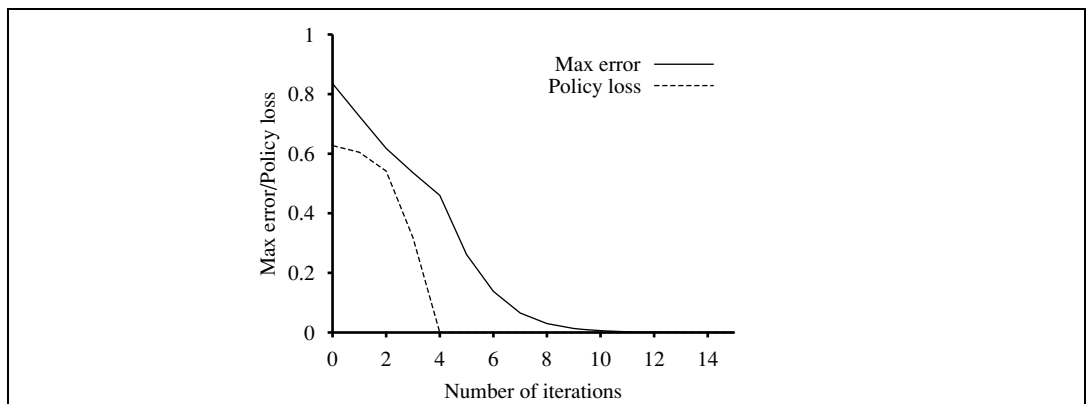
- **Policy evaluation:** given a policy  $\pi_i$ , calculate  $U_i = U^{\pi_i}$ , the utility of each state if  $\pi_i$  were to be executed.

POLICY IMPROVEMENT

- **Policy improvement:** Calculate a new MEU policy  $\pi_{i+1}$ , using one-step look-ahead based on  $U_i$  (as in Equation (17.4)).

The algorithm terminates when the policy improvement step yields no change in the utilities. At this point, we know that the utility function  $U_i$  is a fixed point of the Bellman update, so it is a solution to the Bellman equations, and  $\pi_i$  must be an optimal policy. Because there are only finitely many policies for a finite state space, and each iteration can be shown to yield a better policy, policy iteration must terminate. The algorithm is shown in Figure 17.7.

The policy improvement step is obviously straightforward, but how do we implement the POLICY-EVALUATION routine? It turns out that doing so is much simpler than solving the standard Bellman equations (which is what value iteration does), because the action in each state is fixed by the policy. At the  $i$ th iteration, the policy  $\pi_i$  specifies the action  $\pi_i(s)$  in



**Figure 17.6** The maximum error  $\|U_i - U\|$  of the utility estimates and the policy loss  $\|U^{\pi_i} - U\|$ , as a function of the number of iterations of value iteration.

state  $s$ . This means that we have a simplified version of the Bellman equation (17.5) relating the utility of  $s$  (under  $\pi_i$ ) to the utilities of its neighbors:

$$U_i(s) = R(s) + \gamma \sum_{s'} P(s' | s, \pi_i(s)) U_i(s') . \quad (17.10)$$

For example, suppose  $\pi_i$  is the policy shown in Figure 17.2(a). Then we have  $\pi_i(1, 1) = Up$ ,  $\pi_i(1, 2) = Up$ , and so on, and the simplified Bellman equations are

$$\begin{aligned} U_i(1, 1) &= -0.04 + 0.8U_i(1, 2) + 0.1U_i(1, 1) + 0.1U_i(2, 1) , \\ U_i(1, 2) &= -0.04 + 0.8U_i(1, 3) + 0.2U_i(1, 2) , \\ &\vdots \end{aligned}$$

The important point is that these equations are *linear*, because the “max” operator has been removed. For  $n$  states, we have  $n$  linear equations with  $n$  unknowns, which can be solved exactly in time  $O(n^3)$  by standard linear algebra methods.

For small state spaces, policy evaluation using exact solution methods is often the most efficient approach. For large state spaces,  $O(n^3)$  time might be prohibitive. Fortunately, it is not necessary to do *exact* policy evaluation. Instead, we can perform some number of simplified value iteration steps (simplified because the policy is fixed) to give a reasonably good approximation of the utilities. The simplified Bellman update for this process is

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s' | s, \pi_i(s)) U_i(s') ,$$

and this is repeated  $k$  times to produce the next utility estimate. The resulting algorithm is called **modified policy iteration**. It is often much more efficient than standard policy iteration or value iteration.

MODIFIED POLICY  
ITERATION

```

function POLICY-ITERATION(mdp) returns a policy
  inputs: mdp, an MDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ 
  local variables:  $U$ , a vector of utilities for states in  $S$ , initially zero
                    $\pi$ , a policy vector indexed by state, initially random

  repeat
     $U \leftarrow \text{POLICY-EVALUATION}(\pi, U, \textit{mdp})$ 
     $\textit{unchanged?} \leftarrow \text{true}$ 
    for each state  $s$  in  $S$  do
      if  $\max_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s'] > \sum_{s'} P(s' | s, \pi[s]) U[s']$  then do
         $\pi[s] \leftarrow \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s']$ 
         $\textit{unchanged?} \leftarrow \text{false}$ 
  until  $\textit{unchanged?}$ 
  return  $\pi$ 

```

**Figure 17.7** The policy iteration algorithm for calculating an optimal policy.

The algorithms we have described so far require updating the utility or policy for all states at once. It turns out that this is not strictly necessary. In fact, on each iteration, we can pick *any subset* of states and apply *either* kind of updating (policy improvement or simplified value iteration) to that subset. This very general algorithm is called **asynchronous policy iteration**. Given certain conditions on the initial policy and initial utility function, asynchronous policy iteration is guaranteed to converge to an optimal policy. The freedom to choose any states to work on means that we can design much more efficient heuristic algorithms—for example, algorithms that concentrate on updating the values of states that are likely to be reached by a good policy. This makes a lot of sense in real life: if one has no intention of throwing oneself off a cliff, one should not spend time worrying about the exact value of the resulting states.

## 17.4 PARTIALLY OBSERVABLE MDPs

The description of Markov decision processes in Section 17.1 assumed that the environment was **fully observable**. With this assumption, the agent always knows which state it is in. This, combined with the Markov assumption for the transition model, means that the optimal policy depends only on the current state. When the environment is only **partially observable**, the situation is, one might say, much less clear. The agent does not necessarily know which state it is in, so it cannot execute the action  $\pi(s)$  recommended for that state. Furthermore, the utility of a state  $s$  and the optimal action in  $s$  depend not just on  $s$ , but also on *how much the agent knows* when it is in  $s$ . For these reasons, **partially observable MDPs** (or POMDPs—pronounced “pom-dee-pees”) are usually viewed as much more difficult than ordinary MDPs. We cannot avoid POMDPs, however, because the real world is one.

### 17.4.1 Definition of POMDPs

To get a handle on POMDPs, we must first define them properly. A POMDP has the same elements as an MDP—the transition model  $P(s' | s, a)$ , actions  $A(s)$ , and reward function  $R(s)$ —but, like the partially observable search problems of Section 4.4, it also has a **sensor model**  $P(e | s)$ . Here, as in Chapter 15, the sensor model specifies the probability of perceiving evidence  $e$  in state  $s$ .<sup>3</sup> For example, we can convert the  $4 \times 3$  world of Figure 17.1 into a POMDP by adding a noisy or partial sensor instead of assuming that the agent knows its location exactly. Such a sensor might measure the *number of adjacent walls*, which happens to be 2 in all the nonterminal squares except for those in the third column, where the value is 1; a noisy version might give the wrong value with probability 0.1.

In Chapters 4 and 11, we studied nondeterministic and partially observable planning problems and identified the **belief state**—the set of actual states the agent might be in—as a key concept for describing and calculating solutions. In POMDPs, the belief state  $b$  becomes a *probability distribution* over all possible states, just as in Chapter 15. For example, the initial

<sup>3</sup> As with the reward function for MDPs, the sensor model can also depend on the action and outcome state, but again this change is not fundamental.

belief state for the  $4 \times 3$  POMDP could be the uniform distribution over the nine nonterminal states, i.e.,  $\langle \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, 0, 0 \rangle$ . We write  $b(s)$  for the probability assigned to the actual state  $s$  by belief state  $b$ . The agent can calculate its current belief state as the conditional probability distribution over the actual states given the sequence of percepts and actions so far. This is essentially the **filtering** task described in Chapter 15. The basic recursive filtering equation (15.5 on page 572) shows how to calculate the new belief state from the previous belief state and the new evidence. For POMDPs, we also have an action to consider, but the result is essentially the same. If  $b(s)$  was the previous belief state, and the agent does action  $a$  and then perceives evidence  $e$ , then the new belief state is given by

$$b'(s') = \alpha P(e | s') \sum_s P(s' | s, a) b(s) ,$$

where  $\alpha$  is a normalizing constant that makes the belief state sum to 1. By analogy with the update operator for filtering (page 572), we can write this as

$$b' = \text{FORWARD}(b, a, e) . \quad (17.11)$$

In the  $4 \times 3$  POMDP, suppose the agent moves *Left* and its sensor reports 1 adjacent wall; then it's quite likely (although not guaranteed, because both the motion and the sensor are noisy) that the agent is now in (3,1). Exercise 17.13 asks you to calculate the exact probability values for the new belief state.



The fundamental insight required to understand POMDPs is this: *the optimal action depends only on the agent's current belief state*. That is, the optimal policy can be described by a mapping  $\pi^*(b)$  from belief states to actions. It does *not* depend on the *actual* state the agent is in. This is a good thing, because the agent does not know its actual state; all it knows is the belief state. Hence, the decision cycle of a POMDP agent can be broken down into the following three steps:

1. Given the current belief state  $b$ , execute the action  $a = \pi^*(b)$ .
2. Receive percept  $e$ .
3. Set the current belief state to  $\text{FORWARD}(b, a, e)$  and repeat.

Now we can think of POMDPs as requiring a search in belief-state space, just like the methods for sensorless and contingency problems in Chapter 4. The main difference is that the POMDP belief-state space is *continuous*, because a POMDP belief state is a probability distribution. For example, a belief state for the  $4 \times 3$  world is a point in an 11-dimensional continuous space. An action changes the belief state, not just the physical state. Hence, the action is evaluated at least in part according to the information the agent acquires as a result. POMDPs therefore include the value of information (Section 16.6) as one component of the decision problem.

Let's look more carefully at the outcome of actions. In particular, let's calculate the probability that an agent in belief state  $b$  reaches belief state  $b'$  after executing action  $a$ . Now, if we knew the action *and the subsequent percept*, then Equation (17.11) would provide a *deterministic* update to the belief state:  $b' = \text{FORWARD}(b, a, e)$ . Of course, the subsequent percept is not yet known, so the agent might arrive in one of several possible belief states  $b'$ , depending on the percept that is received. The probability of perceiving  $e$ , given that  $a$  was

performed starting in belief state  $b$ , is given by summing over all the actual states  $s'$  that the agent might reach:

$$\begin{aligned} P(e|a, b) &= \sum_{s'} P(e|a, s', b) P(s'|a, b) \\ &= \sum_{s'} P(e | s') P(s'|a, b) \\ &= \sum_{s'} P(e | s') \sum_s P(s' | s, a) b(s) . \end{aligned}$$

Let us write the probability of reaching  $b'$  from  $b$ , given action  $a$ , as  $P(b' | b, a)$ . Then that gives us

$$\begin{aligned} P(b' | b, a) &= P(b'|a, b) = \sum_e P(b'|e, a, b) P(e|a, b) \\ &= \sum_e P(b'|e, a, b) \sum_{s'} P(e | s') \sum_s P(s' | s, a) b(s) , \end{aligned} \quad (17.12)$$

where  $P(b'|e, a, b)$  is 1 if  $b' = \text{FORWARD}(b, a, e)$  and 0 otherwise.

Equation (17.12) can be viewed as defining a transition model for the belief-state space. We can also define a reward function for belief states (i.e., the expected reward for the actual states the agent might be in):

$$\rho(b) = \sum_s b(s) R(s) .$$

Together,  $P(b' | b, a)$  and  $\rho(b)$  define an *observable* MDP on the space of belief states. Furthermore, it can be shown that an optimal policy for this MDP,  $\pi^*(b)$ , is also an optimal policy for the original POMDP. In other words, *solving a POMDP on a physical state space can be reduced to solving an MDP on the corresponding belief-state space*. This fact is perhaps less surprising if we remember that the belief state is always observable to the agent, by definition.

Notice that, although we have reduced POMDPs to MDPs, the MDP we obtain has a continuous (and usually high-dimensional) state space. None of the MDP algorithms described in Sections 17.2 and 17.3 applies directly to such MDPs. The next two subsections describe a value iteration algorithm designed specifically for POMDPs and an online decision-making algorithm, similar to those developed for games in Chapter 5.

### 17.4.2 Value iteration for POMDPs

Section 17.2 described a value iteration algorithm that computed one utility value for each state. With infinitely many belief states, we need to be more creative. Consider an optimal policy  $\pi^*$  and its application in a specific belief state  $b$ : the policy generates an action, then, for each subsequent percept, the belief state is updated and a new action is generated, and so on. For this specific  $b$ , therefore, the policy is exactly equivalent to a **conditional plan**, as defined in Chapter 4 for nondeterministic and partially observable problems. Instead of thinking about policies, let us think about conditional plans and how the expected utility of executing a fixed conditional plan varies with the initial belief state. We make two observations:





1. Let the utility of executing a *fixed* conditional plan  $p$  starting in physical state  $s$  be  $\alpha_p(s)$ . Then the expected utility of executing  $p$  in belief state  $b$  is just  $\sum_s b(s)\alpha_p(s)$ , or  $b \cdot \alpha_p$  if we think of them both as vectors. Hence, the expected utility of a fixed conditional plan varies *linearly* with  $b$ ; that is, it corresponds to a hyperplane in belief space.
2. At any given belief state  $b$ , the optimal policy will choose to execute the conditional plan with highest expected utility; and the expected utility of  $b$  under the optimal policy is just the utility of that conditional plan:

$$U(b) = U^{\pi^*}(b) = \max_p b \cdot \alpha_p.$$

If the optimal policy  $\pi^*$  chooses to execute  $p$  starting at  $b$ , then it is reasonable to expect that it might choose to execute  $p$  in belief states that are very close to  $b$ ; in fact, if we bound the depth of the conditional plans, then there are only finitely many such plans and the continuous space of belief states will generally be divided into *regions*, each corresponding to a particular conditional plan that is optimal in that region.

From these two observations, we see that the utility function  $U(b)$  on belief states, being the maximum of a collection of hyperplanes, will be *piecewise linear* and *convex*.

To illustrate this, we use a simple two-state world. The states are labeled 0 and 1, with  $R(0) = 0$  and  $R(1) = 1$ . There are two actions: *Stay* stays put with probability 0.9 and *Go* switches to the other state with probability 0.9. For now we will assume the discount factor  $\gamma = 1$ . The sensor reports the correct state with probability 0.6. Obviously, the agent should *Stay* when it thinks it's in state 1 and *Go* when it thinks it's in state 0.

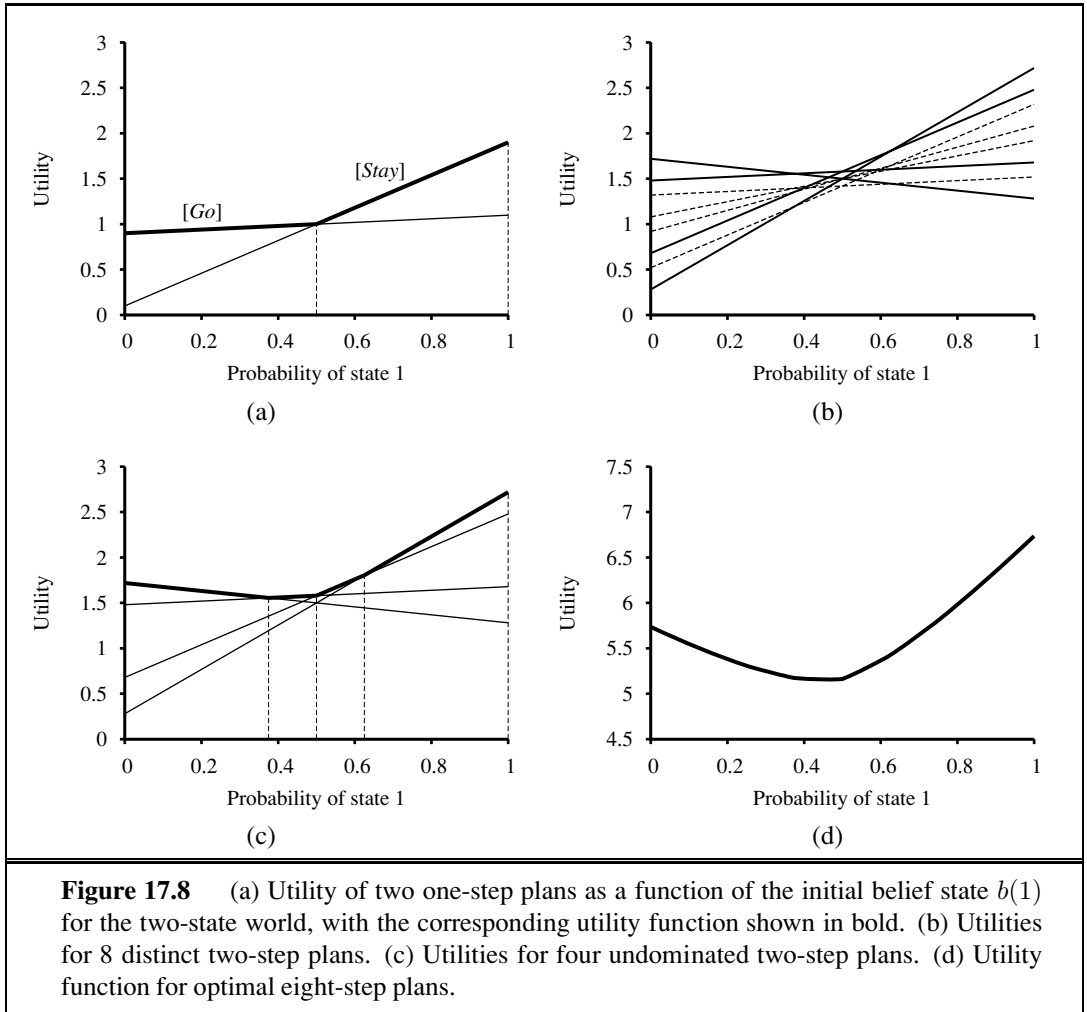
The advantage of a two-state world is that the belief space can be viewed as one-dimensional, because the two probabilities must sum to 1. In Figure 17.8(a), the  $x$ -axis represents the belief state, defined by  $b(1)$ , the probability of being in state 1. Now let us consider the one-step plans  $[Stay]$  and  $[Go]$ , each of which receives the reward for the current state followed by the (discounted) reward for the state reached after the action:

$$\begin{aligned}\alpha_{[Stay]}(0) &= R(0) + \gamma(0.9R(0) + 0.1R(1)) = 0.1 \\ \alpha_{[Stay]}(1) &= R(1) + \gamma(0.9R(1) + 0.1R(0)) = 1.9 \\ \alpha_{[Go]}(0) &= R(0) + \gamma(0.9R(1) + 0.1R(0)) = 0.9 \\ \alpha_{[Go]}(1) &= R(1) + \gamma(0.9R(0) + 0.1R(1)) = 1.1\end{aligned}$$

The hyperplanes (lines, in this case) for  $b \cdot \alpha_{[Stay]}$  and  $b \cdot \alpha_{[Go]}$  are shown in Figure 17.8(a) and their maximum is shown in bold. The bold line therefore represents the utility function for the finite-horizon problem that allows just one action, and in each “piece” of the piecewise linear utility function the optimal action is the first action of the corresponding conditional plan. In this case, the optimal one-step policy is to *Stay* when  $b(1) > 0.5$  and *Go* otherwise.

Once we have utilities  $\alpha_p(s)$  for all the conditional plans  $p$  of depth 1 in each physical state  $s$ , we can compute the utilities for conditional plans of depth 2 by considering each possible first action, each possible subsequent percept, and then each way of choosing a depth-1 plan to execute for each percept:

```
[Stay; if Percept = 0 then Stay else Stay]
[Stay; if Percept = 0 then Stay else Go] ...
```



There are eight distinct depth-2 plans in all, and their utilities are shown in Figure 17.8(b). Notice that four of the plans, shown as dashed lines, are suboptimal across the entire belief space—we say these plans are **dominated**, and they need not be considered further. There are four undominated plans, each of which is optimal in a specific region, as shown in Figure 17.8(c). The regions partition the belief-state space.

We repeat the process for depth 3, and so on. In general, let  $p$  be a depth- $d$  conditional plan whose initial action is  $a$  and whose depth- $d - 1$  subplan for percept  $e$  is  $p.e$ ; then

$$\alpha_p(s) = R(s) + \gamma \left( \sum_{s'} P(s' | s, a) \sum_e P(e | s') \alpha_{p.e}(s') \right). \quad (17.13)$$

This recursion naturally gives us a value iteration algorithm, which is sketched in Figure 17.9. The structure of the algorithm and its error analysis are similar to those of the basic value iteration algorithm in Figure 17.4 on page 653; the main difference is that instead of computing one utility number for each state, POMDP-VALUE-ITERATION maintains a collection of

```

function POMDP-VALUE-ITERATION(pomdp,  $\epsilon$ ) returns a utility function
  inputs: pomdp, a POMDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ ,
           sensor model  $P(e | s)$ , rewards  $R(s)$ , discount  $\gamma$ 
            $\epsilon$ , the maximum error allowed in the utility of any state
  local variables:  $U, U'$ , sets of plans  $p$  with associated utility vectors  $\alpha_p$ 

   $U' \leftarrow$  a set containing just the empty plan  $[\ ]$ , with  $\alpha_{[\ ]}(s) = R(s)$ 
  repeat
     $U \leftarrow U'$ 
     $U' \leftarrow$  the set of all plans consisting of an action and, for each possible next percept,
                a plan in  $U$  with utility vectors computed according to Equation (17.13)
     $U' \leftarrow$  REMOVE-DOMINATED-PLANS( $U'$ )
  until MAX-DIFFERENCE( $U, U'$ )  $< \epsilon(1 - \gamma)/\gamma$ 
  return  $U$ 

```

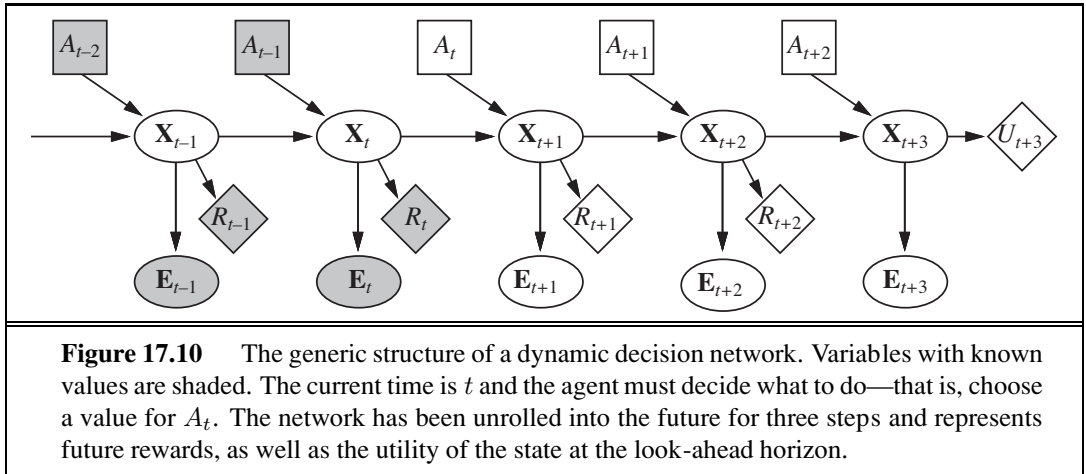
**Figure 17.9** A high-level sketch of the value iteration algorithm for POMDPs. The REMOVE-DOMINATED-PLANS step and MAX-DIFFERENCE test are typically implemented as linear programs.

undominated plans with their utility hyperplanes. The algorithm's complexity depends primarily on how many plans get generated. Given  $|A|$  actions and  $|E|$  possible observations, it is easy to show that there are  $|A|^{O(|E|^{d-1})}$  distinct depth- $d$  plans. Even for the lowly two-state world with  $d = 8$ , the exact number is  $2^{255}$ . The elimination of dominated plans is essential for reducing this doubly exponential growth: the number of undominated plans with  $d = 8$  is just 144. The utility function for these 144 plans is shown in Figure 17.8(d).

Notice that even though state 0 has lower utility than state 1, the intermediate belief states have even lower utility because the agent lacks the information needed to choose a good action. This is why information has value in the sense defined in Section 16.6 and optimal policies in POMDPs often include information-gathering actions.

Given such a utility function, an executable policy can be extracted by looking at which hyperplane is optimal at any given belief state  $b$  and executing the first action of the corresponding plan. In Figure 17.8(d), the corresponding optimal policy is still the same as for depth-1 plans: *Stay* when  $b(1) > 0.5$  and *Go* otherwise.

In practice, the value iteration algorithm in Figure 17.9 is hopelessly inefficient for larger problems—even the  $4 \times 3$  POMDP is too hard. The main reason is that, given  $n$  conditional plans at level  $d$ , the algorithm constructs  $|A| \cdot n^{|E|}$  conditional plans at level  $d + 1$  before eliminating the dominated ones. Since the 1970s, when this algorithm was developed, there have been several advances including more efficient forms of value iteration and various kinds of policy iteration algorithms. Some of these are discussed in the notes at the end of the chapter. For general POMDPs, however, finding optimal policies is very difficult (PSPACE-hard, in fact—i.e., very hard indeed). Problems with a few dozen states are often infeasible. The next section describes a different, approximate method for solving POMDPs, one based on look-ahead search.



### 17.4.3 Online agents for POMDPs

In this section, we outline a simple approach to agent design for partially observable, stochastic environments. The basic elements of the design are already familiar:

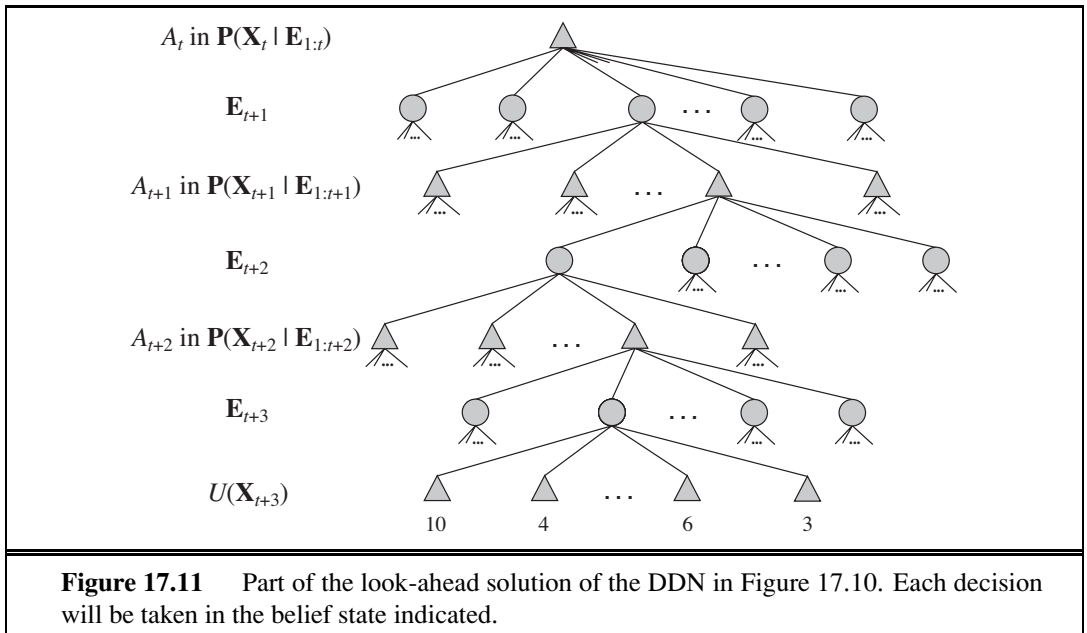
- The transition and sensor models are represented by a **dynamic Bayesian network** (DBN), as described in Chapter 15.
- The dynamic Bayesian network is extended with decision and utility nodes, as used in **decision networks** in Chapter 16. The resulting model is called a **dynamic decision network**, or DDN.
- A filtering algorithm is used to incorporate each new percept and action and to update the belief state representation.
- Decisions are made by projecting forward possible action sequences and choosing the best one.

DYNAMIC DECISION  
NETWORK

DBNs are **factored representations** in the terminology of Chapter 2; they typically have an exponential complexity advantage over atomic representations and can model quite substantial real-world problems. The agent design is therefore a practical implementation of the **utility-based agent** sketched in Chapter 2.

In the DBN, the single state  $S_t$  becomes a set of state variables  $\mathbf{X}_t$ , and there may be multiple evidence variables  $\mathbf{E}_t$ . We will use  $A_t$  to refer to the action at time  $t$ , so the transition model becomes  $\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{X}_t, A_t)$  and the sensor model becomes  $\mathbf{P}(\mathbf{E}_t|\mathbf{X}_t)$ . We will use  $R_t$  to refer to the reward received at time  $t$  and  $U_t$  to refer to the utility of the state at time  $t$ . (Both of these are random variables.) With this notation, a dynamic decision network looks like the one shown in Figure 17.10.

Dynamic decision networks can be used as inputs for any POMDP algorithm, including those for value and policy iteration methods. In this section, we focus on look-ahead methods that project action sequences forward from the current belief state in much the same way as do the game-playing algorithms of Chapter 5. The network in Figure 17.10 has been projected three steps into the future; the current and future decisions  $A$  and the future observations



$\mathbf{E}$  and rewards  $R$  are all unknown. Notice that the network includes nodes for the *rewards* for  $\mathbf{X}_{t+1}$  and  $\mathbf{X}_{t+2}$ , but the *utility* for  $\mathbf{X}_{t+3}$ . This is because the agent must maximize the (discounted) sum of all future rewards, and  $U(\mathbf{X}_{t+3})$  represents the reward for  $\mathbf{X}_{t+3}$  and all subsequent rewards. As in Chapter 5, we assume that  $U$  is available only in some approximate form: if exact utility values were available, look-ahead beyond depth 1 would be unnecessary.

Figure 17.11 shows part of the search tree corresponding to the three-step look-ahead DDN in Figure 17.10. Each of the triangular nodes is a belief state in which the agent makes a decision  $A_{t+i}$  for  $i = 0, 1, 2, \dots$ . The round (chance) nodes correspond to choices by the environment, namely, what evidence  $\mathbf{E}_{t+i}$  arrives. Notice that there are no chance nodes corresponding to the action outcomes; this is because the belief-state update for an action is deterministic regardless of the actual outcome.

The belief state at each triangular node can be computed by applying a filtering algorithm to the sequence of percepts and actions leading to it. In this way, the algorithm takes into account the fact that, for decision  $A_{t+i}$ , the agent *will* have available percepts  $\mathbf{E}_{t+1}, \dots, \mathbf{E}_{t+i}$ , even though at time  $t$  it does not know what those percepts will be. In this way, a decision-theoretic agent automatically takes into account the value of information and will execute information-gathering actions where appropriate.

A decision can be extracted from the search tree by backing up the utility values from the leaves, taking an average at the chance nodes and taking the maximum at the decision nodes. This is similar to the EXPECTIMINIMAX algorithm for game trees with chance nodes, except that (1) there can also be rewards at non-leaf states and (2) the decision nodes correspond to belief states rather than actual states. The time complexity of an exhaustive search to depth  $d$  is  $O(|A|^d \cdot |\mathbf{E}|^d)$ , where  $|A|$  is the number of available actions and  $|\mathbf{E}|$  is the number of possible percepts. (Notice that this is far less than the number of depth- $d$  conditional

plans generated by value iteration.) For problems in which the discount factor  $\gamma$  is not too close to 1, a shallow search is often good enough to give near-optimal decisions. It is also possible to approximate the averaging step at the chance nodes, by sampling from the set of possible percepts instead of summing over all possible percepts. There are various other ways of finding good approximate solutions quickly, but we defer them to Chapter 21.

Decision-theoretic agents based on dynamic decision networks have a number of advantages compared with other, simpler agent designs presented in earlier chapters. In particular, they handle partially observable, uncertain environments and can easily revise their “plans” to handle unexpected evidence. With appropriate sensor models, they can handle sensor failure and can plan to gather information. They exhibit “graceful degradation” under time pressure and in complex environments, using various approximation techniques. So what is missing? One defect of our DDN-based algorithm is its reliance on forward search through state space, rather than using the hierarchical and other advanced planning techniques described in Chapter 11. There have been attempts to extend these techniques into the probabilistic domain, but so far they have proved to be inefficient. A second, related problem is the basically propositional nature of the DDN language. We would like to be able to extend some of the ideas for first-order probabilistic languages to the problem of decision making. Current research has shown that this extension is possible and has significant benefits, as discussed in the notes at the end of the chapter.

## 17.5 DECISIONS WITH MULTIPLE AGENTS: GAME THEORY

### GAME THEORY

This chapter has concentrated on making decisions in uncertain environments. But what if the uncertainty is due to other agents and the decisions they make? And what if the decisions of those agents are in turn influenced by our decisions? We addressed this question once before, when we studied games in Chapter 5. There, however, we were primarily concerned with turn-taking games in fully observable environments, for which minimax search can be used to find optimal moves. In this section we study the aspects of **game theory** that analyze games with simultaneous moves and other sources of partial observability. (Game theorists use the terms **perfect information** and **imperfect information** rather than fully and partially observable.) Game theory can be used in at least two ways:

1. **Agent design:** Game theory can analyze the agent’s decisions and compute the expected utility for each decision (under the assumption that other agents are acting optimally according to game theory). For example, in the game **two-finger Morra**, two players,  $O$  and  $E$ , simultaneously display one or two fingers. Let the total number of fingers be  $f$ . If  $f$  is odd,  $O$  collects  $f$  dollars from  $E$ ; and if  $f$  is even,  $E$  collects  $f$  dollars from  $O$ . Game theory can determine the best strategy against a rational player and the expected return for each player.<sup>4</sup>

<sup>4</sup> Morra is a recreational version of an **inspection game**. In such games, an inspector chooses a day to inspect a facility (such as a restaurant or a biological weapons plant), and the facility operator chooses a day to hide all the nasty stuff. The inspector wins if the days are different, and the facility operator wins if they are the same.

2. **Mechanism design:** When an environment is inhabited by many agents, it might be possible to define the rules of the environment (i.e., the game that the agents must play) so that the collective good of all agents is maximized when each agent adopts the game-theoretic solution that maximizes its own utility. For example, game theory can help design the protocols for a collection of Internet traffic routers so that each router has an incentive to act in such a way that global throughput is maximized. Mechanism design can also be used to construct intelligent **multiagent systems** that solve complex problems in a distributed fashion.

17.5.1 Single-move games

We start by considering a restricted set of games: ones where all players take action simultaneously and the result of the game is based on this single set of actions. (Actually, it is not crucial that the actions take place at exactly the same time; what matters is that no player has knowledge of the other players’ choices.) The restriction to a single move (and the very use of the word “game”) might make this seem trivial, but in fact, game theory is serious business. It is used in decision-making situations including the auctioning of oil drilling rights and wireless frequency spectrum rights, bankruptcy proceedings, product development and pricing decisions, and national defense—situations involving billions of dollars and hundreds of thousands of lives. A single-move game is defined by three components:

- PLAYER
- ACTION
- PAYOFF FUNCTION
- STRATEGIC FORM
- **Players** or agents who will be making decisions. Two-player games have received the most attention, although  $n$ -player games for  $n > 2$  are also common. We give players capitalized names, like *Alice* and *Bob* or *O* and *E*.
  - **Actions** that the players can choose. We will give actions lowercase names, like *one* or *testify*. The players may or may not have the same set of actions available.
  - **A payoff function** that gives the utility to each player for each combination of actions by all the players. For single-move games the payoff function can be represented by a matrix, a representation known as the **strategic form** (also called **normal form**). The payoff matrix for two-finger Morra is as follows:

	<i>O: one</i>	<i>O: two</i>
<i>E: one</i>	$E = +2, O = -2$	$E = -3, O = +3$
<i>E: two</i>	$E = -3, O = +3$	$E = +4, O = -4$

For example, the lower-right corner shows that when player *O* chooses action *two* and *E* also chooses *two*, the payoff is +4 for *E* and −4 for *O*.

- STRATEGY
- PURE STRATEGY
- MIXED STRATEGY
- STRATEGY PROFILE
- OUTCOME
- Each player in a game must adopt and then execute a **strategy** (which is the name used in game theory for a *policy*). A **pure strategy** is a deterministic policy; for a single-move game, a pure strategy is just a single action. For many games an agent can do better with a **mixed strategy**, which is a randomized policy that selects actions according to a probability distribution. The mixed strategy that chooses action  $a$  with probability  $p$  and action  $b$  otherwise is written  $[p:a; (1 - p):b]$ . For example, a mixed strategy for two-finger Morra might be  $[0.5:one; 0.5:two]$ . A **strategy profile** is an assignment of a strategy to each player; given the strategy profile, the game’s **outcome** is a numeric value for each player.

SOLUTION

A **solution** to a game is a strategy profile in which each player adopts a rational strategy. We will see that the most important issue in game theory is to define what “rational” means when each agent chooses only part of the strategy profile that determines the outcome. It is important to realize that outcomes are actual results of playing a game, while solutions are theoretical constructs used to analyze a game. We will see that some games have a solution only in mixed strategies. But that does not mean that a player must literally be adopting a mixed strategy to be rational.

PRISONER'S  
DILEMMA

Consider the following story: Two alleged burglars, Alice and Bob, are caught red-handed near the scene of a burglary and are interrogated separately. A prosecutor offers each a deal: if you testify against your partner as the leader of a burglary ring, you’ll go free for being the cooperative one, while your partner will serve 10 years in prison. However, if you both testify against each other, you’ll both get 5 years. Alice and Bob also know that if both refuse to testify they will serve only 1 year each for the lesser charge of possessing stolen property. Now Alice and Bob face the so-called **prisoner’s dilemma**: should they testify or refuse? Being rational agents, Alice and Bob each want to maximize their own expected utility. Let’s assume that Alice is callously unconcerned about her partner’s fate, so her utility decreases in proportion to the number of years she will spend in prison, regardless of what happens to Bob. Bob feels exactly the same way. To help reach a rational decision, they both construct the following payoff matrix:

	<i>Alice: testify</i>	<i>Alice: refuse</i>
<i>Bob: testify</i>	$A = -5, B = -5$	$A = -10, B = 0$
<i>Bob: refuse</i>	$A = 0, B = -10$	$A = -1, B = -1$

DOMINANT  
STRATEGY  
STRONG  
DOMINATION

Alice analyzes the payoff matrix as follows: “Suppose Bob testifies. Then I get 5 years if I testify and 10 years if I don’t, so in that case testifying is better. On the other hand, if Bob refuses, then I get 0 years if I testify and 1 year if I refuse, so in that case as well testifying is better. So in either case, it’s better for me to testify, so that’s what I must do.”

WEAK DOMINATION

Alice has discovered that *testify* is a **dominant strategy** for the game. We say that a strategy  $s$  for player  $p$  **strongly dominates** strategy  $s'$  if the outcome for  $s$  is better for  $p$  than the outcome for  $s'$ , for every choice of strategies by the other player(s). Strategy  $s$  **weakly dominates**  $s'$  if  $s$  is better than  $s'$  on at least one strategy profile and no worse on any other. A dominant strategy is a strategy that dominates all others. It is irrational to play a dominated strategy, and irrational not to play a dominant strategy if one exists. Being rational, Alice chooses the dominant strategy. We need just a bit more terminology: we say that an outcome is **Pareto optimal**<sup>5</sup> if there is no other outcome that all players would prefer. An outcome is **Pareto dominated** by another outcome if all players would prefer the other outcome.

PERETO OPTIMAL  
PERETO DOMINATED

If Alice is clever as well as rational, she will continue to reason as follows: Bob’s dominant strategy is also to testify. Therefore, he will testify and we will both get five years. When each player has a dominant strategy, the combination of those strategies is called a **dominant strategy equilibrium**. In general, a strategy profile forms an **equilibrium** if no player can benefit by switching strategies, given that every other player sticks with the same

DOMINANT  
STRATEGY  
EQUILIBRIUM  
EQUILIBRIUM

<sup>5</sup> Pareto optimality is named after the economist Vilfredo Pareto (1848–1923).





## NASH EQUILIBRIUM

strategy. An equilibrium is essentially a **local optimum** in the space of policies; it is the top of a peak that slopes downward along every dimension, where a dimension corresponds to a player's strategy choices.

The mathematician John Nash (1928–) proved that *every game has at least one equilibrium*. The general concept of equilibrium is now called **Nash equilibrium** in his honor. Clearly, a dominant strategy equilibrium is a Nash equilibrium (Exercise 17.16), but some games have Nash equilibria but no dominant strategies.

The *dilemma* in the prisoner's dilemma is that the equilibrium outcome is worse for both players than the outcome they would get if they both refused to testify. In other words,  $(testify, testify)$  is Pareto dominated by the  $(-1, -1)$  outcome of  $(refuse, refuse)$ . Is there any way for Alice and Bob to arrive at the  $(-1, -1)$  outcome? It is certainly an *allowable* option for both of them to refuse to testify, but it is hard to see how rational agents can get there, given the definition of the game. Either player contemplating playing *refuse* will realize that he or she would do better by playing *testify*. That is the attractive power of an equilibrium point. Game theorists agree that being a Nash equilibrium is a necessary condition for being a solution—although they disagree whether it is a sufficient condition.

It is easy enough to get to the  $(refuse, refuse)$  solution if we modify the game. For example, we could change to a **repeated game** in which the players know that they will meet again. Or the agents might have moral beliefs that encourage cooperation and fairness. That means they have a different utility function, necessitating a different payoff matrix, making it a different game. We will see later that agents with limited computational powers, rather than the ability to reason absolutely rationally, can reach non-equilibrium outcomes, as can an agent that knows that the other agent has limited rationality. In each case, we are considering a different game than the one described by the payoff matrix above.

Now let's look at a game that has no dominant strategy. Acme, a video game console manufacturer, has to decide whether its next game machine will use Blu-ray discs or DVDs. Meanwhile, the video game software producer Best needs to decide whether to produce its next game on Blu-ray or DVD. The profits for both will be positive if they agree and negative if they disagree, as shown in the following payoff matrix:

	<i>Acme:bluray</i>	<i>Acme:dvd</i>
<i>Best:bluray</i>	$A = +9, B = +9$	$A = -4, B = -1$
<i>Best:dvd</i>	$A = -3, B = -1$	$A = +5, B = +5$

There is no dominant strategy equilibrium for this game, but there are *two* Nash equilibria:  $(bluray, bluray)$  and  $(dvd, dvd)$ . We know these are Nash equilibria because if either player unilaterally moves to a different strategy, that player will be worse off. Now the agents have a problem: *there are multiple acceptable solutions, but if each agent aims for a different solution, then both agents will suffer*. How can they agree on a solution? One answer is that both should choose the Pareto-optimal solution  $(bluray, bluray)$ ; that is, we can restrict the definition of “solution” to the unique Pareto-optimal Nash equilibrium *provided that one exists*. Every game has at least one Pareto-optimal solution, but a game might have several, or they might not be equilibrium points. For example, if  $(bluray, bluray)$  had payoff  $(5, 5)$ , then there would be two equal Pareto-optimal equilibrium points. To choose between



COORDINATION  
GAME

them the agents can either guess or *communicate*, which can be done either by establishing a convention that orders the solutions before the game begins or by negotiating to reach a mutually beneficial solution during the game (which would mean including communicative actions as part of a sequential game). Communication thus arises in game theory for exactly the same reasons that it arose in multiagent planning in Section 11.4. Games in which players need to communicate like this are called **coordination games**.

A game can have more than one Nash equilibrium; how do we know that every game must have at least one? Some games have no *pure-strategy* Nash equilibria. Consider, for example, any pure-strategy profile for two-finger Morra (page 666). If the total number of fingers is even, then *O* will want to switch; on the other hand (so to speak), if the total is odd, then *E* will want to switch. Therefore, no pure strategy profile can be an equilibrium and we must look to mixed strategies instead.

ZERO-SUM GAME

But *which* mixed strategy? In 1928, von Neumann developed a method for finding the *optimal* mixed strategy for two-player, **zero-sum games**—games in which the sum of the payoffs is always zero.<sup>6</sup> Clearly, Morra is such a game. For two-player, zero-sum games, we know that the payoffs are equal and opposite, so we need consider the payoffs of only one player, who will be the maximizer (just as in Chapter 5). For Morra, we pick the even player *E* to be the maximizer, so we can define the payoff matrix by the values  $U_E(e, o)$ —the payoff to *E* if *E* does *e* and *O* does *o*. (For convenience we call player *E* “her” and *O* “him.”) Von Neumann’s method is called the **maximin** technique, and it works as follows:

MAXIMIN

- Suppose we change the rules as follows: first *E* picks her strategy and reveals it to *O*. Then *O* picks his strategy, with knowledge of *E*’s strategy. Finally, we evaluate the expected payoff of the game based on the chosen strategies. This gives us a turn-taking game to which we can apply the standard **minimax** algorithm from Chapter 5. Let’s suppose this gives an outcome  $U_{E,O}$ . Clearly, this game favors *O*, so the true utility  $U$  of the original game (from *E*’s point of view) is *at least*  $U_{E,O}$ . For example, if we just look at pure strategies, the minimax game tree has a root value of  $-3$  (see Figure 17.12(a)), so we know that  $U \geq -3$ .
- Now suppose we change the rules to force *O* to reveal his strategy first, followed by *E*. Then the minimax value of this game is  $U_{O,E}$ , and because this game favors *E* we know that  $U$  is *at most*  $U_{O,E}$ . With pure strategies, the value is  $+2$  (see Figure 17.12(b)), so we know  $U \leq +2$ .

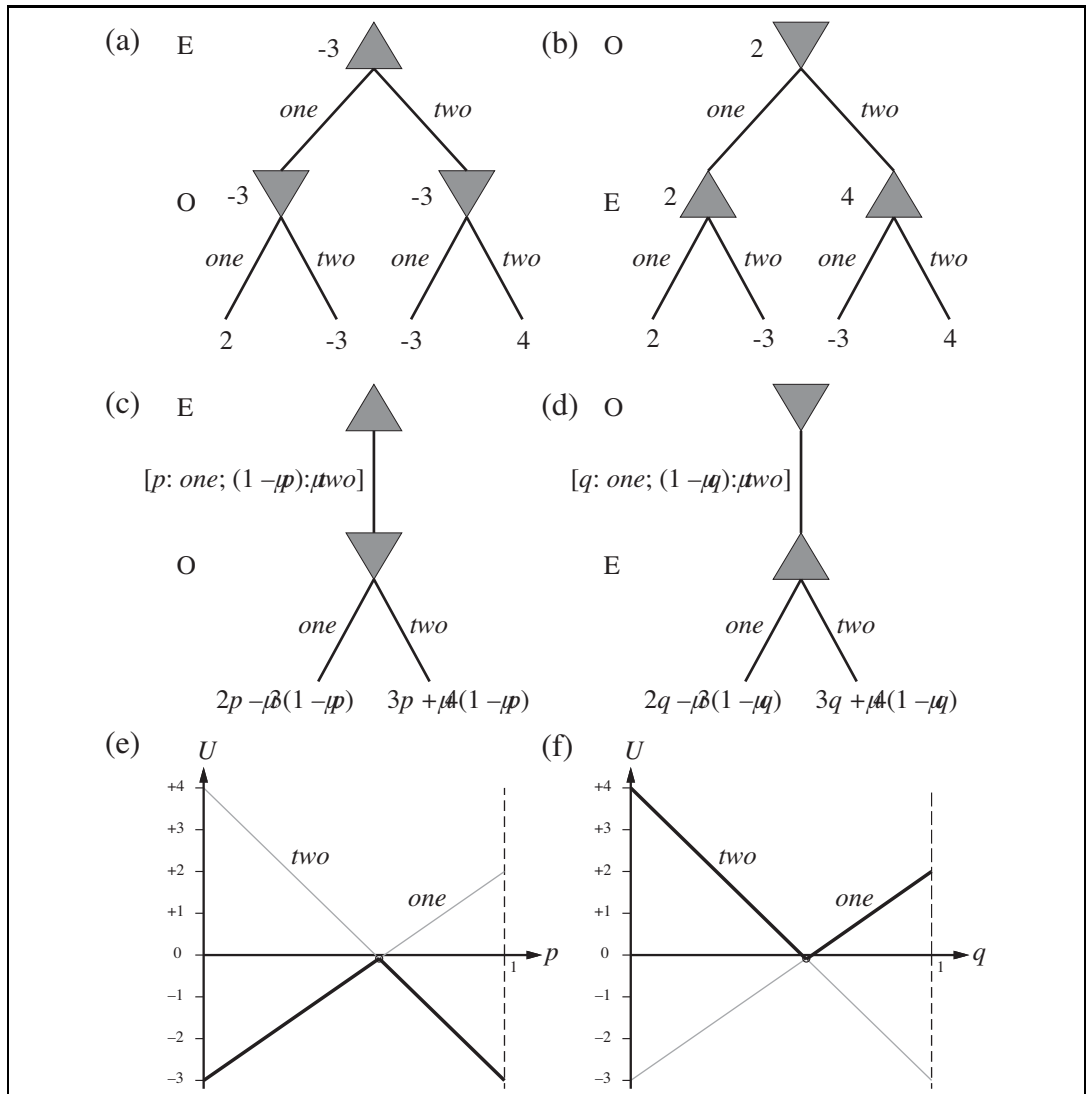
Combining these two arguments, we see that the true utility  $U$  of the solution to the original game must satisfy

$$U_{E,O} \leq U \leq U_{O,E} \quad \text{or in this case,} \quad -3 \leq U \leq 2.$$



To pinpoint the value of  $U$ , we need to turn our analysis to mixed strategies. First, observe the following: *once the first player has revealed his or her strategy, the second player might as well choose a pure strategy*. The reason is simple: if the second player plays a mixed strategy,  $[p: \text{one}; (1-p): \text{two}]$ , its expected utility is a linear combination  $(p \cdot u_{\text{one}} + (1-p) \cdot u_{\text{two}})$  of

<sup>6</sup> or a constant—see page 162.



**Figure 17.12** (a) and (b): Minimax game trees for two-finger Morra if the players take turns playing pure strategies. (c) and (d): Parameterized game trees where the first player plays a mixed strategy. The payoffs depend on the probability parameter ( $p$  or  $q$ ) in the mixed strategy. (e) and (f): For any particular value of the probability parameter, the second player will choose the “better” of the two actions, so the value of the first player’s mixed strategy is given by the heavy lines. The first player will choose the probability parameter for the mixed strategy at the intersection point.

the utilities of the pure strategies,  $u_{\text{one}}$  and  $u_{\text{two}}$ . This linear combination can never be better than the better of  $u_{\text{one}}$  and  $u_{\text{two}}$ , so the second player can just choose the better one.

With this observation in mind, the minimax trees can be thought of as having infinitely many branches at the root, corresponding to the infinitely many mixed strategies the first

player can choose. Each of these leads to a node with two branches corresponding to the pure strategies for the second player. We can depict these infinite trees finitely by having one “parameterized” choice at the root:

- If  $E$  chooses first, the situation is as shown in Figure 17.12(c).  $E$  chooses the strategy  $[p: \text{one}; (1-p): \text{two}]$  at the root, and then  $O$  chooses a pure strategy (and hence a move) given the value of  $p$ . If  $O$  chooses *one*, the expected payoff (to  $E$ ) is  $2p - 3(1-p) = 5p - 3$ ; if  $O$  chooses *two*, the expected payoff is  $-3p + 4(1-p) = 4 - 7p$ . We can draw these two payoffs as straight lines on a graph, where  $p$  ranges from 0 to 1 on the  $x$ -axis, as shown in Figure 17.12(e).  $O$ , the minimizer, will always choose the lower of the two lines, as shown by the heavy lines in the figure. Therefore, the best that  $E$  can do at the root is to choose  $p$  to be at the intersection point, which is where

$$5p - 3 = 4 - 7p \quad \Rightarrow \quad p = 7/12.$$

The utility for  $E$  at this point is  $U_{E,O} = -1/12$ .

- If  $O$  moves first, the situation is as shown in Figure 17.12(d).  $O$  chooses the strategy  $[q: \text{one}; (1-q): \text{two}]$  at the root, and then  $E$  chooses a move given the value of  $q$ . The payoffs are  $2q - 3(1-q) = 5q - 3$  and  $-3q + 4(1-q) = 4 - 7q$ .<sup>7</sup> Again, Figure 17.12(f) shows that the best  $O$  can do at the root is to choose the intersection point:

$$5q - 3 = 4 - 7q \quad \Rightarrow \quad q = 7/12.$$

The utility for  $E$  at this point is  $U_{O,E} = -1/12$ .

Now we know that the true utility of the original game lies between  $-1/12$  and  $-1/12$ , that is, it is exactly  $-1/12$ ! (The moral is that it is better to be  $O$  than  $E$  if you are playing this game.) Furthermore, the true utility is attained by the mixed strategy  $[7/12: \text{one}; 5/12: \text{two}]$ , which should be played by both players. This strategy is called the **maximin equilibrium** of the game, and is a Nash equilibrium. Note that each component strategy in an equilibrium mixed strategy has the same expected utility. In this case, both *one* and *two* have the same expected utility,  $-1/12$ , as the mixed strategy itself.

Our result for two-finger Morra is an example of the general result by von Neumann: *every two-player zero-sum game has a maximin equilibrium when you allow mixed strategies*. Furthermore, every Nash equilibrium in a zero-sum game is a maximin for both players. A player who adopts the maximin strategy has two guarantees: First, no other strategy can do better against an opponent who plays well (although some other strategies might be better at exploiting an opponent who makes irrational mistakes). Second, the player continues to do just as well even if the strategy is revealed to the opponent.

The general algorithm for finding maximin equilibria in zero-sum games is somewhat more involved than Figures 17.12(e) and (f) might suggest. When there are  $n$  possible actions, a mixed strategy is a point in  $n$ -dimensional space and the lines become hyperplanes. It's also possible for some pure strategies for the second player to be dominated by others, so that they are not optimal against *any* strategy for the first player. After removing all such strategies (which might have to be done repeatedly), the optimal choice at the root is the

<sup>7</sup> It is a coincidence that these equations are the same as those for  $p$ ; the coincidence arises because  $U_E(\text{one}, \text{two}) = U_E(\text{two}, \text{one}) = -3$ . This also explains why the optimal strategy is the same for both players.



highest (or lowest) intersection point of the remaining hyperplanes. Finding this choice is an example of a **linear programming** problem: maximizing an objective function subject to linear constraints. Such problems can be solved by standard techniques in time polynomial in the number of actions (and in the number of bits used to specify the reward function, if you want to get technical).

The question remains, what should a rational agent actually *do* in playing a single game of Morra? The rational agent will have derived the fact that  $[7/12: one; 5/12: two]$  is the maximin equilibrium strategy, and will assume that this is mutual knowledge with a rational opponent. The agent could use a 12-sided die or a random number generator to pick randomly according to this mixed strategy, in which case the expected payoff would be  $-1/12$  for *E*. Or the agent could just decide to play *one*, or *two*. In either case, the expected payoff remains  $-1/12$  for *E*. Curiously, unilaterally choosing a particular action does not harm one's expected payoff, but allowing the other agent to know that one has made such a unilateral decision *does* affect the expected payoff, because then the opponent can adjust his strategy accordingly.

Finding equilibria in non-zero-sum games is somewhat more complicated. The general approach has two steps: (1) Enumerate all possible subsets of actions that might form mixed strategies. For example, first try all strategy profiles where each player uses a single action, then those where each player uses either one or two actions, and so on. This is exponential in the number of actions, and so only applies to relatively small games. (2) For each strategy profile enumerated in (1), check to see if it is an equilibrium. This is done by solving a set of equations and inequalities that are similar to the ones used in the zero-sum case. For two players these equations are linear and can be solved with basic linear programming techniques, but for three or more players they are nonlinear and may be very difficult to solve.

## 17.5.2 Repeated games

### REPEATED GAME

So far we have looked only at games that last a single move. The simplest kind of multiple-move game is the **repeated game**, in which players face the same choice repeatedly, but each time with knowledge of the history of all players' previous choices. A strategy profile for a repeated game specifies an action choice for each player at each time step for every possible history of previous choices. As with MDPs, payoffs are additive over time.

Let's consider the repeated version of the prisoner's dilemma. Will Alice and Bob work together and refuse to testify, knowing they will meet again? The answer depends on the details of the engagement. For example, suppose Alice and Bob know that they must play exactly 100 rounds of prisoner's dilemma. Then they both know that the 100th round will not be a repeated game—that is, its outcome can have no effect on future rounds—and therefore they will both choose the dominant strategy, *testify*, in that round. But once the 100th round is determined, the 99th round can have no effect on subsequent rounds, so it too will have a dominant strategy equilibrium at (*testify*, *testify*). By induction, both players will choose *testify* on every round, earning a total jail sentence of 500 years each.

We can get different solutions by changing the rules of the interaction. For example, suppose that after each round there is a 99% chance that the players will meet again. Then the expected number of rounds is still 100, but neither player knows for sure which round

PERPETUAL  
PUNISHMENT

will be the last. Under these conditions, more cooperative behavior is possible. For example, one equilibrium strategy is for each player to *refuse* unless the other player has ever played *testify*. This strategy could be called **perpetual punishment**. Suppose both players have adopted this strategy, and this is mutual knowledge. Then as long as neither player has played *testify*, then at any point in time the expected future total payoff for each player is

$$\sum_{t=0}^{\infty} 0.99^t \cdot (-1) = -100.$$

A player who deviates from the strategy and chooses *testify* will gain a score of 0 rather than  $-1$  on the very next move, but from then on both players will play *testify* and the player's total expected future payoff becomes

$$0 + \sum_{t=1}^{\infty} 0.99^t \cdot (-5) = -495.$$

Therefore, at every step, there is no incentive to deviate from (*refuse*, *refuse*). Perpetual punishment is the “mutually assured destruction” strategy of the prisoner's dilemma: once either player decides to *testify*, it ensures that both players suffer a great deal. But it works as a deterrent only if the other player believes you have adopted this strategy—or at least that you might have adopted it.

TIT-FOR-TAT

Other strategies are more forgiving. The most famous, called **tit-for-tat**, calls for starting with *refuse* and then echoing the other player's previous move on all subsequent moves. So Alice would refuse as long as Bob refuses and would testify the move after Bob testified, but would go back to refusing if Bob did. Although very simple, this strategy has proven to be highly robust and effective against a wide variety of strategies.

We can also get different solutions by changing the agents, rather than changing the rules of engagement. Suppose the agents are finite-state machines with  $n$  states and they are playing a game with  $m > n$  total steps. The agents are thus incapable of representing the number of remaining steps, and must treat it as an unknown. Therefore, they cannot do the induction, and are free to arrive at the more favorable (*refuse*, *refuse*) equilibrium. In this case, ignorance is bliss—or rather, having your opponent believe that you are ignorant is bliss. Your success in these repeated games depends on the other player's *perception* of you as a bully or a simpleton, and not on your actual characteristics.

### 17.5.3 Sequential games

EXTENSIVE FORM

In the general case, a game consists of a sequence of turns that need not be all the same. Such games are best represented by a game tree, which game theorists call the **extensive form**. The tree includes all the same information we saw in Section 5.1: an initial state  $S_0$ , a function  $\text{PLAYER}(s)$  that tells which player has the move, a function  $\text{ACTIONS}(s)$  enumerating the possible actions, a function  $\text{RESULT}(s, a)$  that defines the transition to a new state, and a partial function  $\text{UTILITY}(s, p)$ , which is defined only on terminal states, to give the payoff for each player.

To represent stochastic games, such as backgammon, we add a distinguished player, *chance*, that can take random actions. *Chance*'s “strategy” is part of the definition of the

game, specified as a probability distribution over actions (the other players get to choose their own strategy). To represent games with nondeterministic actions, such as billiards, we break the action into two pieces: the player's action itself has a deterministic result, and then *chance* has a turn to react to the action in its own capricious way. To represent simultaneous moves, as in the prisoner's dilemma or two-finger Morra, we impose an arbitrary order on the players, but we have the option of asserting that the earlier player's actions are not observable to the subsequent players: e.g., Alice must choose *refuse* or *testify* first, then Bob chooses, but Bob does not know what choice Alice made at that time (we can also represent the fact that the move is revealed later). However, we assume the players always remember all their *own* previous actions; this assumption is called **perfect recall**.

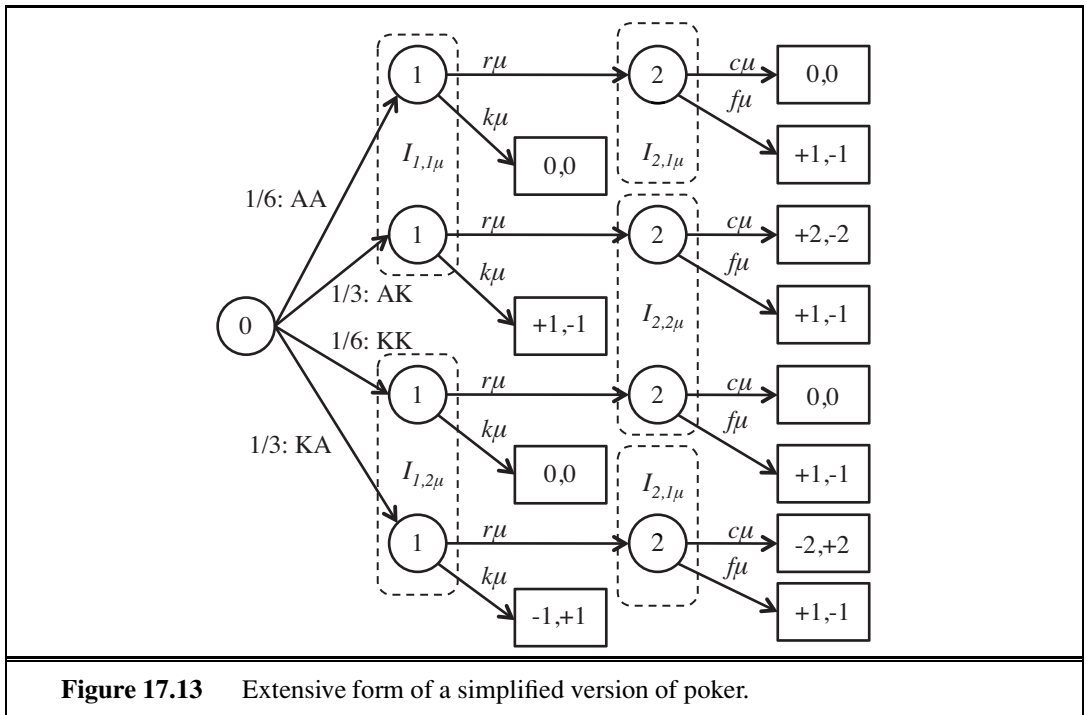
The key idea of extensive form that sets it apart from the game trees of Chapter 5 is the representation of partial observability. We saw in Section 5.6 that a player in a partially observable game such as Kriegspiel can create a game tree over the space of **belief states**. With that tree, we saw that in some cases a player can find a sequence of moves (a strategy) that leads to a forced checkmate regardless of what actual state we started in, and regardless of what strategy the opponent uses. However, the techniques of Chapter 5 could not tell a player what to do when there is no guaranteed checkmate. If the player's best strategy depends on the opponent's strategy and vice versa, then minimax (or alpha-beta) by itself cannot find a solution. The extensive form *does* allow us to find solutions because it represents the belief states (game theorists call them **information sets**) of *all* players at once. From that representation we can find equilibrium solutions, just as we did with normal-form games.

INFORMATION SETS

As a simple example of a sequential game, place two agents in the  $4 \times 3$  world of Figure 17.1 and have them move simultaneously until one agent reaches an exit square, and gets the payoff for that square. If we specify that no movement occurs when the two agents try to move into the same square simultaneously (a common problem at many traffic intersections), then certain pure strategies can get stuck forever. Thus, agents need a mixed strategy to perform well in this game: randomly choose between moving ahead and staying put. This is exactly what is done to resolve packet collisions in Ethernet networks.

Next we'll consider a very simple variant of poker. The deck has only four cards, two aces and two kings. One card is dealt to each player. The first player then has the option to *raise* the stakes of the game from 1 point to 2, or to *check*. If player 1 checks, the game is over. If he raises, then player 2 has the option to *call*, accepting that the game is worth 2 points, or *fold*, conceding the 1 point. If the game does not end with a fold, then the payoff depends on the cards: it is zero for both players if they have the same card; otherwise the player with the king pays the stakes to the player with the ace.

The extensive-form tree for this game is shown in Figure 17.13. Nonterminal states are shown as circles, with the player to move inside the circle; player 0 is *chance*. Each action is depicted as an arrow with a label, corresponding to a *raise*, *check*, *call*, or *fold*, or, for *chance*, the four possible deals ("AK" means that player 1 gets an ace and player 2 a king). Terminal states are rectangles labeled by their payoff to player 1 and player 2. Information sets are shown as labeled dashed boxes; for example,  $I_{1,1}$  is the information set where it is player 1's turn, and he knows he has an ace (but does not know what player 2 has). In information set  $I_{2,1}$ , it is player 2's turn and she knows that she has an ace and that player 1 has raised,



but does not know what card player 1 has. (Due to the limits of two-dimensional paper, this information set is shown as two boxes rather than one.)

One way to solve an extensive game is to convert it to a normal-form game. Recall that the normal form is a matrix, each row of which is labeled with a pure strategy for player 1, and each column by a pure strategy for player 2. In an extensive game a pure strategy for player  $i$  corresponds to an action for each information set involving that player. So in Figure 17.13, one pure strategy for player 1 is “raise when in  $I_{1,1}$  (that is, when I have an ace), and check when in  $I_{1,2}$  (when I have a king).” In the payoff matrix below, this strategy is called  $rk$ . Similarly, strategy  $cf$  for player 2 means “call when I have an ace and fold when I have a king.” Since this is a zero-sum game, the matrix below gives only the payoff for player 1; player 2 always has the opposite payoff:

	2:cc	2:cf	2:ff	2:fc
1:rr	0	-1/6	1	7/6
1:kr	-1/3	-1/6	5/6	2/3
1:rk	1/3	<b>0</b>	1/6	1/2
1:kk	0	<b>0</b>	0	0

This game is so simple that it has two pure-strategy equilibria, shown in bold:  $cf$  for player 2 and  $rk$  or  $kk$  for player 1. But in general we can solve extensive games by converting to normal form and then finding a solution (usually a mixed strategy) using standard linear programming methods. That works in theory. But if a player has  $I$  information sets and  $a$  actions per set, then that player will have  $a^I$  pure strategies. In other words, the size of the normal-form matrix is exponential in the number of information sets, so in practice the



approach works only for very small game trees, on the order of a dozen states. A game like Texas hold'em poker has about  $10^{18}$  states, making this approach completely infeasible.

What are the alternatives? In Chapter 5 we saw how alpha-beta search could handle games of perfect information with huge game trees by generating the tree incrementally, by pruning some branches, and by heuristically evaluating nonterminal nodes. But that approach does not work well for games with imperfect information, for two reasons: first, it is harder to prune, because we need to consider mixed strategies that combine multiple branches, not a pure strategy that always chooses the best branch. Second, it is harder to heuristically evaluate a nonterminal node, because we are dealing with information sets, not individual states.

SEQUENCE FORM

Koller *et al.* (1996) come to the rescue with an alternative representation of extensive games, called the **sequence form**, that is only linear in the size of the tree, rather than exponential. Rather than represent strategies, it represents paths through the tree; the number of paths is equal to the number of terminal nodes. Standard linear programming methods can again be applied to this representation. The resulting system can solve poker variants with 25,000 states in a minute or two. This is an exponential speedup over the normal-form approach, but still falls far short of handling full poker, with  $10^{18}$  states.

ABSTRACTION

If we can't handle  $10^{18}$  states, perhaps we can simplify the problem by changing the game to a simpler form. For example, if I hold an ace and am considering the possibility that the next card will give me a pair of aces, then I don't care about the suit of the next card; any suit will do equally well. This suggests forming an **abstraction** of the game, one in which suits are ignored. The resulting game tree will be smaller by a factor of  $4! = 24$ . Suppose I can solve this smaller game; how will the solution to that game relate to the original game? If no player is going for a flush (or bluffing so), then the suits don't matter to any player, and the solution for the abstraction will also be a solution for the original game. However, if any player is contemplating a flush, then the abstraction will be only an approximate solution (but it is possible to compute bounds on the error).

There are many opportunities for abstraction. For example, at the point in a game where each player has two cards, if I hold a pair of queens, then the other players' hands could be abstracted into three classes: *better* (only a pair of kings or a pair of aces), *same* (pair of queens) or *worse* (everything else). However, this abstraction might be too coarse. A better abstraction would divide *worse* into, say, *medium pair* (nines through jacks), *low pair*, and *no pair*. These examples are abstractions of states; it is also possible to abstract actions. For example, instead of having a bet action for each integer from 1 to 1000, we could restrict the bets to  $10^0$ ,  $10^1$ ,  $10^2$  and  $10^3$ . Or we could cut out one of the rounds of betting altogether. We can also abstract over chance nodes, by considering only a subset of the possible deals. This is equivalent to the rollout technique used in Go programs. Putting all these abstractions together, we can reduce the  $10^{18}$  states of poker to  $10^7$  states, a size that can be solved with current techniques.

Poker programs based on this approach can easily defeat novice and some experienced human players, but are not yet at the level of master players. Part of the problem is that the solution these programs approximate—the equilibrium solution—is optimal only against an opponent who also plays the equilibrium strategy. Against fallible human players it is important to be able to exploit an opponent's deviation from the equilibrium strategy. As

Gautam Rao (aka “The Count”), the world’s leading online poker player, said (Billings *et al.*, 2003), “You have a very strong program. Once you add opponent modeling to it, it will kill everyone.” However, good models of human fallability remain elusive.

In a sense, extensive game form is the one of the most complete representations we have seen so far: it can handle partially observable, multiagent, stochastic, sequential, dynamic environments—most of the hard cases from the list of environment properties on page 42. However, there are two limitations of game theory. First, it does not deal well with continuous states and actions (although there have been some extensions to the continuous case; for example, the theory of **Cournot competition** uses game theory to solve problems where two companies choose prices for their products from a continuous space). Second, game theory assumes the game is *known*. Parts of the game may be specified as unobservable to some of the players, but it must be known what parts are unobservable. In cases in which the players learn the unknown structure of the game over time, the model begins to break down. Let’s examine each source of uncertainty, and whether each can be represented in game theory.

**Actions:** There is no easy way to represent a game where the players have to discover what actions are available. Consider the game between computer virus writers and security experts. Part of the problem is anticipating what action the virus writers will try next.

**Strategies:** Game theory is very good at representing the idea that the other players’ strategies are initially unknown—as long as we assume all agents are rational. The theory itself does not say what to do when the other players are less than fully rational. The notion of a **Bayes–Nash equilibrium** partially addresses this point: it is an equilibrium with respect to a player’s prior probability distribution over the other players’ strategies—in other words, it expresses a player’s beliefs about the other players’ likely strategies.

**Chance:** If a game depends on the roll of a die, it is easy enough to model a chance node with uniform distribution over the outcomes. But what if it is possible that the die is unfair? We can represent that with another chance node, higher up in the tree, with two branches for “die is fair” and “die is unfair,” such that the corresponding nodes in each branch are in the same information set (that is, the players don’t know if the die is fair or not). And what if we suspect the other opponent does know? Then we add *another* chance node, with one branch representing the case where the opponent does know, and one where he doesn’t.

**Utilities:** What if we don’t know our opponent’s utilities? Again, that can be modeled with a chance node, such that the other agent knows its own utilities in each branch, but we don’t. But what if we don’t know our *own* utilities? For example, how do I know if it is rational to order the Chef’s salad if I don’t know how much I will like it? We can model that with yet another chance node specifying an unobservable “intrinsic quality” of the salad.

Thus, we see that game theory is good at representing most sources of uncertainty—but at the cost of doubling the size of the tree every time we add another node; a habit which quickly leads to intractably large trees. Because of these and other problems, game theory has been used primarily to *analyze* environments that are at equilibrium, rather than to *control* agents within an environment. Next we shall see how it can help *design* environments.

COURNOT  
COMPETITION

BAYES–NASH  
EQUILIBRIUM

## 17.6 MECHANISM DESIGN

MECHANISM DESIGN

In the previous section, we asked, “Given a game, what is a rational strategy?” In this section, we ask, “Given that agents pick rational strategies, what game should we design?” More specifically, we would like to design a game whose solutions, consisting of each agent pursuing its own rational strategy, result in the maximization of some global utility function. This problem is called **mechanism design**, or sometimes **inverse game theory**. Mechanism design is a staple of economics and political science. Capitalism 101 says that if everyone tries to get rich, the total wealth of society will increase. But the examples we will discuss show that proper mechanism design is necessary to keep the invisible hand on track. For collections of agents, mechanism design allows us to construct smart systems out of a collection of more limited systems—even uncooperative systems—in much the same way that teams of humans can achieve goals beyond the reach of any individual.

MECHANISM  
CENTER

Examples of mechanism design include auctioning off cheap airline tickets, routing TCP packets between computers, deciding how medical interns will be assigned to hospitals, and deciding how robotic soccer players will cooperate with their teammates. Mechanism design became more than an academic subject in the 1990s when several nations, faced with the problem of auctioning off licenses to broadcast in various frequency bands, lost hundreds of millions of dollars in potential revenue as a result of poor mechanism design. Formally, a **mechanism** consists of (1) a language for describing the set of allowable strategies that agents may adopt, (2) a distinguished agent, called the **center**, that collects reports of strategy choices from the agents in the game, and (3) an outcome rule, known to all agents, that the center uses to determine the payoffs to each agent, given their strategy choices.

### 17.6.1 Auctions

AUCTION

Let’s consider **auctions** first. An auction is a mechanism for selling some goods to members of a pool of bidders. For simplicity, we concentrate on auctions with a single item for sale. Each bidder  $i$  has a utility value  $v_i$  for having the item. In some cases, each bidder has a **private value** for the item. For example, the first item sold on eBay was a broken laser pointer, which sold for \$14.83 to a collector of broken laser pointers. Thus, we know that the collector has  $v_i \geq \$14.83$ , but most other people would have  $v_j \ll \$14.83$ . In other cases, such as auctioning drilling rights for an oil tract, the item has a **common value**—the tract will produce some amount of money,  $X$ , and all bidders value a dollar equally—but there is uncertainty as to what the actual value of  $X$  is. Different bidders have different information, and hence different estimates of the item’s true value. In either case, bidders end up with their own  $v_i$ . Given  $v_i$ , each bidder gets a chance, at the appropriate time or times in the auction, to make a bid  $b_i$ . The highest bid,  $b_{max}$  wins the item, but the price paid need not be  $b_{max}$ ; that’s part of the mechanism design.

ASCENDING-BID  
ENGLISH AUCTION

The best-known auction mechanism is the **ascending-bid**,<sup>8</sup> or **English auction**, in which the center starts by asking for a minimum (or **reserve**) bid  $b_{min}$ . If some bidder is

<sup>8</sup> The word “auction” comes from the Latin *augere*, to increase.

willing to pay that amount, the center then asks for  $b_{min} + d$ , for some increment  $d$ , and continues up from there. The auction ends when nobody is willing to bid anymore; then the last bidder wins the item, paying the price he bid.

EFFICIENT

How do we know if this is a good mechanism? One goal is to maximize expected revenue for the seller. Another goal is to maximize a notion of global utility. These goals overlap to some extent, because one aspect of maximizing global utility is to ensure that the winner of the auction is the agent who values the item the most (and thus is willing to pay the most). We say an auction is **efficient** if the goods go to the agent who values them most. The ascending-bid auction is usually both efficient and revenue maximizing, but if the reserve price is set too high, the bidder who values it most may not bid, and if the reserve is set too low, the seller loses net revenue.

COLLUSION

Probably the most important things that an auction mechanism can do is encourage a sufficient number of bidders to enter the game and discourage them from engaging in **collusion**. Collusion is an unfair or illegal agreement by two or more bidders to manipulate prices. It can happen in secret backroom deals or tacitly, within the rules of the mechanism.

For example, in 1999, Germany auctioned ten blocks of cell-phone spectrum with a simultaneous auction (bids were taken on all ten blocks at the same time), using the rule that any bid must be a minimum of a 10% raise over the previous bid on a block. There were only two credible bidders, and the first, Mannesman, entered the bid of 20 million deutschmark on blocks 1-5 and 18.18 million on blocks 6-10. Why 18.18M? One of T-Mobile's managers said they "interpreted Mannesman's first bid as an offer." Both parties could compute that a 10% raise on 18.18M is 19.99M; thus Mannesman's bid was interpreted as saying "we can each get half the blocks for 20M; let's not spoil it by bidding the prices up higher." And in fact T-Mobile bid 20M on blocks 6-10 and that was the end of the bidding. The German government got less than they expected, because the two competitors were able to use the bidding mechanism to come to a tacit agreement on how not to compete. From the government's point of view, a better result could have been obtained by any of these changes to the mechanism: a higher reserve price; a sealed-bid first-price auction, so that the competitors could not communicate through their bids; or incentives to bring in a third bidder. Perhaps the 10% rule was an error in mechanism design, because it facilitated the precise signaling from Mannesman to T-Mobile.

STRATEGY-PROOF

TRUTH-REVEALING

REVELATION  
PRINCIPLE

In general, both the seller and the global utility function benefit if there are more bidders, although global utility can suffer if you count the cost of wasted time of bidders that have no chance of winning. One way to encourage more bidders is to make the mechanism easier for them. After all, if it requires too much research or computation on the part of the bidders, they may decide to take their money elsewhere. So it is desirable that the bidders have a **dominant strategy**. Recall that "dominant" means that the strategy works against all other strategies, which in turn means that an agent can adopt it without regard for the other strategies. An agent with a dominant strategy can just bid, without wasting time contemplating other agents' possible strategies. A mechanism where agents have a dominant strategy is called a **strategy-proof** mechanism. If, as is usually the case, that strategy involves the bidders revealing their true value,  $v_i$ , then it is called a **truth-revealing**, or **truthful**, auction; the term **incentive compatible** is also used. The **revelation principle** states that any mecha-

nism can be transformed into an equivalent truth-revealing mechanism, so part of mechanism design is finding these equivalent mechanisms.

It turns out that the ascending-bid auction has most of the desirable properties. The bidder with the highest value  $v_i$  gets the goods at a price of  $b_o + d$ , where  $b_o$  is the highest bid among all the other agents and  $d$  is the auctioneer's increment.<sup>9</sup> Bidders have a simple dominant strategy: keep bidding as long as the current cost is below your  $v_i$ . The mechanism is not quite truth-revealing, because the winning bidder reveals only that his  $v_i \geq b_o + d$ ; we have a lower bound on  $v_i$  but not an exact amount.

A disadvantage (from the point of view of the seller) of the ascending-bid auction is that it can discourage competition. Suppose that in a bid for cell-phone spectrum there is one advantaged company that everyone agrees would be able to leverage existing customers and infrastructure, and thus can make a larger profit than anyone else. Potential competitors can see that they have no chance in an ascending-bid auction, because the advantaged company can always bid higher. Thus, the competitors may not enter at all, and the advantaged company ends up winning at the reserve price.

Another negative property of the English auction is its high communication costs. Either the auction takes place in one room or all bidders have to have high-speed, secure communication lines; in either case they have to have the time available to go through several rounds of bidding. An alternative mechanism, which requires much less communication, is the **sealed-bid auction**. Each bidder makes a single bid and communicates it to the auctioneer, without the other bidders seeing it. With this mechanism, there is no longer a simple dominant strategy. If your value is  $v_i$  and you believe that the maximum of all the other agents' bids will be  $b_o$ , then you should bid  $b_o + \epsilon$ , for some small  $\epsilon$ , if that is less than  $v_i$ . Thus, your bid depends on your estimation of the other agents' bids, requiring you to do more work. Also, note that the agent with the highest  $v_i$  might not win the auction. This is offset by the fact that the auction is more competitive, reducing the bias toward an advantaged bidder.

A small change in the mechanism for sealed-bid auctions produces the **sealed-bid second-price auction**, also known as a **Vickrey auction**.<sup>10</sup> In such auctions, the winner pays the price of the *second*-highest bid,  $b_o$ , rather than paying his own bid. This simple modification completely eliminates the complex deliberations required for standard (or **first-price**) sealed-bid auctions, because the dominant strategy is now simply to bid  $v_i$ ; the mechanism is truth-revealing. Note that the utility of agent  $i$  in terms of his bid  $b_i$ , his value  $v_i$ , and the best bid among the other agents,  $b_o$ , is

$$u_i = \begin{cases} (v_i - b_o) & \text{if } b_i > b_o \\ 0 & \text{otherwise.} \end{cases}$$

To see that  $b_i = v_i$  is a dominant strategy, note that when  $(v_i - b_o)$  is positive, any bid that wins the auction is optimal, and bidding  $v_i$  in particular wins the auction. On the other hand, when  $(v_i - b_o)$  is negative, any bid that loses the auction is optimal, and bidding  $v_i$  in

<sup>9</sup> There is actually a small chance that the agent with highest  $v_i$  fails to get the goods, in the case in which  $b_o < v_i < b_o + d$ . The chance of this can be made arbitrarily small by decreasing the increment  $d$ .

<sup>10</sup> Named after William Vickrey (1914–1996), who won the 1996 Nobel Prize in economics for this work and died of a heart attack three days later.

SEALED-BID  
AUCTION

SEALED-BID  
SECOND-PRICE  
AUCTION  
VICKREY AUCTION

particular loses the auction. So bidding  $v_i$  is optimal for all possible values of  $b_o$ , and in fact,  $v_i$  is the only bid that has this property. Because of its simplicity and the minimal computation requirements for both seller and bidders, the Vickrey auction is widely used in constructing distributed AI systems. Also, Internet search engines conduct over a billion auctions a day to sell advertisements along with their search results, and online auction sites handle \$100 billion a year in goods, all using variants of the Vickrey auction. Note that the expected value to the seller is  $b_o$ , which is the same expected return as the limit of the English auction as the increment  $d$  goes to zero. This is actually a very general result: the **revenue equivalence theorem** states that, with a few minor caveats, any auction mechanism where risk-neutral bidders have values  $v_i$  known only to themselves (but know a probability distribution from which those values are sampled), will yield the same expected revenue. This principle means that the various mechanisms are not competing on the basis of revenue generation, but rather on other qualities.

Although the second-price auction is truth-revealing, it turns out that extending the idea to multiple goods and using a next-price auction is not truth-revealing. Many Internet search engines use a mechanism where they auction  $k$  slots for ads on a page. The highest bidder wins the top spot, the second highest gets the second spot, and so on. Each winner pays the price bid by the next-lower bidder, with the understanding that payment is made only if the searcher actually clicks on the ad. The top slots are considered more valuable because they are more likely to be noticed and clicked on. Imagine that three bidders,  $b_1$ ,  $b_2$  and  $b_3$ , have valuations for a click of  $v_1 = 200$ ,  $v_2 = 180$ , and  $v_3 = 100$ , and that  $k = 2$  slots are available, where it is known that the top spot is clicked on 5% of the time and the bottom spot 2%. If all bidders bid truthfully, then  $b_1$  wins the top slot and pays 180, and has an expected return of  $(200 - 180) \times 0.05 = 1$ . The second slot goes to  $b_2$ . But  $b_1$  can see that if she were to bid anything in the range 101–179, she would concede the top slot to  $b_2$ , win the second slot, and yield an expected return of  $(200 - 100) \times .02 = 2$ . Thus,  $b_1$  can double her expected return by bidding less than her true value in this case. In general, bidders in this multislot auction must spend a lot of energy analyzing the bids of others to determine their best strategy; there is no simple dominant strategy. Aggarwal *et al.* (2006) show that there is a unique truthful auction mechanism for this multislot problem, in which the winner of slot  $j$  pays the full price for slot  $j$  just for those additional clicks that are available at slot  $j$  and not at slot  $j + 1$ . The winner pays the price for the lower slot for the remaining clicks. In our example,  $b_1$  would bid 200 truthfully, and would pay 180 for the additional  $.05 - .02 = .03$  clicks in the top slot, but would pay only the cost of the bottom slot, 100, for the remaining  $.02$  clicks. Thus, the total return to  $b_1$  would be  $(200 - 180) \times .03 + (200 - 100) \times .02 = 2.6$ .

Another example of where auctions can come into play within AI is when a collection of agents are deciding whether to cooperate on a joint plan. Hunsberger and Grosz (2000) show that this can be accomplished efficiently with an auction in which the agents bid for roles in the joint plan.

### 17.6.2 Common goods

TRAGEDY OF THE  
COMMONS

Now let's consider another type of game, in which countries set their policy for controlling air pollution. Each country has a choice: they can reduce pollution at a cost of -10 points for implementing the necessary changes, or they can continue to pollute, which gives them a net utility of -5 (in added health costs, etc.) and also contributes -1 points to every other country (because the air is shared across countries). Clearly, the dominant strategy for each country is "continue to pollute," but if there are 100 countries and each follows this policy, then each country gets a total utility of -104, whereas if every country reduced pollution, they would each have a utility of -10. This situation is called the **tragedy of the commons**: if nobody has to pay for using a common resource, then it tends to be exploited in a way that leads to a lower total utility for all agents. It is similar to the prisoner's dilemma: there is another solution to the game that is better for all parties, but there appears to be no way for rational agents to arrive at that solution.

EXTERNALITIES

The standard approach for dealing with the tragedy of the commons is to change the mechanism to one that charges each agent for using the commons. More generally, we need to ensure that all **externalities**—effects on global utility that are not recognized in the individual agents' transactions—are made explicit. Setting the prices correctly is the difficult part. In the limit, this approach amounts to creating a mechanism in which each agent is effectively required to maximize global utility, but can do so by making a local decision. For this example, a carbon tax would be an example of a mechanism that charges for use of the commons in a way that, if implemented well, maximizes global utility.

VICKREY-CLARKE-  
GROVES  
VCG

As a final example, consider the problem of allocating some common goods. Suppose a city decides it wants to install some free wireless Internet transceivers. However, the number of transceivers they can afford is less than the number of neighborhoods that want them. The city wants to allocate the goods efficiently, to the neighborhoods that would value them the most. That is, they want to maximize the global utility  $V = \sum_i v_i$ . The problem is that if they just ask each neighborhood council "how much do you value this free gift?" they would all have an incentive to lie, and report a high value. It turns out there is a mechanism, known as the **Vickrey-Clarke-Groves**, or **VCG**, mechanism, that makes it a dominant strategy for each agent to report its true utility and that achieves an efficient allocation of the goods. The trick is that each agent pays a tax equivalent to the loss in global utility that occurs because of the agent's presence in the game. The mechanism works like this:

1. The center asks each agent to report its value for receiving an item. Call this  $b_i$ .
2. The center allocates the goods to a subset of the bidders. We call this subset  $A$ , and use the notation  $b_i(A)$  to mean the result to  $i$  under this allocation:  $b_i$  if  $i$  is in  $A$  (that is,  $i$  is a winner), and 0 otherwise. The center chooses  $A$  to maximize total reported utility  $B = \sum_i b_i(A)$ .
3. The center calculates (for each  $i$ ) the sum of the reported utilities for all the winners except  $i$ . We use the notation  $B_{-i} = \sum_{j \neq i} b_j(A)$ . The center also computes (for each  $i$ ) the allocation that would maximize total global utility if  $i$  were not in the game; call that sum  $W_{-i}$ .
4. Each agent  $i$  pays a tax equal to  $W_{-i} - B_{-i}$ .

In this example, the VCG rule means that each winner would pay a tax equal to the highest reported value among the losers. That is, if I report my value as 5, and that causes someone with value 2 to miss out on an allocation, then I pay a tax of 2. All winners should be happy because they pay a tax that is less than their value, and all losers are as happy as they can be, because they value the goods less than the required tax.

Why is it that this mechanism is truth-revealing? First, consider the payoff to agent  $i$ , which is the value of getting an item, minus the tax:

$$v_i(A) - (W_{-i} - B_{-i}) . \quad (17.14)$$

Here we distinguish the agent's true utility,  $v_i$ , from his reported utility  $b_i$  (but we are trying to show that a dominant strategy is  $b_i = v_i$ ). Agent  $i$  knows that the center will maximize global utility using the reported values,

$$\sum_j b_j(A) = b_i(A) + \sum_{j \neq i} b_j(A)$$

whereas agent  $i$  wants the center to maximize (17.14), which can be rewritten as

$$v_i(A) + \sum_{j \neq i} b_j(A) - W_{-i} .$$

Since agent  $i$  cannot affect the value of  $W_{-i}$  (it depends only on the other agents), the only way  $i$  can make the center optimize what  $i$  wants is to report the true utility,  $b_i = v_i$ .

## 17.7 SUMMARY

---

This chapter shows how to use knowledge about the world to make decisions even when the outcomes of an action are uncertain and the rewards for acting might not be reaped until many actions have passed. The main points are as follows:

- Sequential decision problems in uncertain environments, also called **Markov decision processes**, or MDPs, are defined by a **transition model** specifying the probabilistic outcomes of actions and a **reward function** specifying the reward in each state.
- The utility of a state sequence is the sum of all the rewards over the sequence, possibly discounted over time. The solution of an MDP is a **policy** that associates a decision with every state that the agent might reach. An optimal policy maximizes the utility of the state sequences encountered when it is executed.
- The utility of a state is the expected utility of the state sequences encountered when an optimal policy is executed, starting in that state. The **value iteration** algorithm for solving MDPs works by iteratively solving the equations relating the utility of each state to those of its neighbors.
- **Policy iteration** alternates between calculating the utilities of states under the current policy and improving the current policy with respect to the current utilities.
- Partially observable MDPs, or POMDPs, are much more difficult to solve than are MDPs. They can be solved by conversion to an MDP in the continuous space of belief



states; both value iteration and policy iteration algorithms have been devised. Optimal behavior in POMDPs includes information gathering to reduce uncertainty and therefore make better decisions in the future.

- A decision-theoretic agent can be constructed for POMDP environments. The agent uses a **dynamic decision network** to represent the transition and sensor models, to update its belief state, and to project forward possible action sequences.
- **Game theory** describes rational behavior for agents in situations in which multiple agents interact simultaneously. Solutions of games are **Nash equilibria**—strategy profiles in which no agent has an incentive to deviate from the specified strategy.
- **Mechanism design** can be used to set the rules by which agents will interact, in order to maximize some global utility through the operation of individually rational agents. Sometimes, mechanisms exist that achieve this goal without requiring each agent to consider the choices made by other agents.

We shall return to the world of MDPs and POMDP in Chapter 21, when we study **reinforcement learning** methods that allow an agent to improve its behavior from experience in sequential, uncertain environments.

---

## BIBLIOGRAPHICAL AND HISTORICAL NOTES

Richard Bellman developed the ideas underlying the modern approach to sequential decision problems while working at the RAND Corporation beginning in 1949. According to his autobiography (Bellman, 1984), he coined the exciting term “dynamic programming” to hide from a research-phobic Secretary of Defense, Charles Wilson, the fact that his group was doing mathematics. (This cannot be strictly true, because his first paper using the term (Bellman, 1952) appeared before Wilson became Secretary of Defense in 1953.) Bellman’s book, *Dynamic Programming* (1957), gave the new field a solid foundation and introduced the basic algorithmic approaches. Ron Howard’s Ph.D. thesis (1960) introduced policy iteration and the idea of average reward for solving infinite-horizon problems. Several additional results were introduced by Bellman and Dreyfus (1962). Modified policy iteration is due to van Nunen (1976) and Puterman and Shin (1978). Asynchronous policy iteration was analyzed by Williams and Baird (1993), who also proved the policy loss bound in Equation (17.9). The analysis of discounting in terms of stationary preferences is due to Koopmans (1972). The texts by Bertsekas (1987), Puterman (1994), and Bertsekas and Tsitsiklis (1996) provide a rigorous introduction to sequential decision problems. Papadimitriou and Tsitsiklis (1987) describe results on the computational complexity of MDPs.

Seminal work by Sutton (1988) and Watkins (1989) on reinforcement learning methods for solving MDPs played a significant role in introducing MDPs into the AI community, as did the later survey by Barto *et al.* (1995). (Earlier work by Werbos (1977) contained many similar ideas, but was not taken up to the same extent.) The connection between MDPs and AI planning problems was made first by Sven Koenig (1991), who showed how probabilistic STRIPS operators provide a compact representation for transition models (see also Wellman,

1990b). Work by Dean *et al.* (1993) and Tash and Russell (1994) attempted to overcome the combinatorics of large state spaces by using a limited search horizon and abstract states. Heuristics based on the value of information can be used to select areas of the state space where a local expansion of the horizon will yield a significant improvement in decision quality. Agents using this approach can tailor their effort to handle time pressure and generate some interesting behaviors such as using familiar “beaten paths” to find their way around the state space quickly without having to recompute optimal decisions at each point.

As one might expect, AI researchers have pushed MDPs in the direction of more expressive representations that can accommodate much larger problems than the traditional atomic representations based on transition matrices. The use of a dynamic Bayesian network to represent transition models was an obvious idea, but work on **factored MDPs** (Boutilier *et al.*, 2000; Koller and Parr, 2000; Guestrin *et al.*, 2003b) extends the idea to structured representations of the value function with provable improvements in complexity. **Relational MDPs** (Boutilier *et al.*, 2001; Guestrin *et al.*, 2003a) go one step further, using structured representations to handle domains with many related objects.

The observation that a partially observable MDP can be transformed into a regular MDP over belief states is due to Astrom (1965) and Aoki (1965). The first complete algorithm for the exact solution of POMDPs—essentially the value iteration algorithm presented in this chapter—was proposed by Edward Sondik (1971) in his Ph.D. thesis. (A later journal paper by Smallwood and Sondik (1973) contains some errors, but is more accessible.) Lovejoy (1991) surveyed the first twenty-five years of POMDP research, reaching somewhat pessimistic conclusions about the feasibility of solving large problems. The first significant contribution within AI was the Witness algorithm (Cassandra *et al.*, 1994; Kaelbling *et al.*, 1998), an improved version of POMDP value iteration. Other algorithms soon followed, including an approach due to Hansen (1998) that constructs a policy incrementally in the form of a finite-state automaton. In this policy representation, the belief state corresponds directly to a particular state in the automaton. More recent work in AI has focused on **point-based** value iteration methods that, at each iteration, generate conditional plans and  $\alpha$ -vectors for a finite set of belief states rather than for the entire belief space. Lovejoy (1991) proposed such an algorithm for a fixed grid of points, an approach taken also by Bonet (2002). An influential paper by Pineau *et al.* (2003) suggested generating reachable points by simulating trajectories in a somewhat greedy fashion; Spaan and Vlassis (2005) observe that one need generate plans for only a small, randomly selected subset of points to improve on the plans from the previous iteration for all points in the set. Current point-based methods—such as point-based policy iteration (Ji *et al.*, 2007)—can generate near-optimal solutions for POMDPs with thousands of states. Because POMDPs are PSPACE-hard (Papadimitriou and Tsitsiklis, 1987), further progress may require taking advantage of various kinds of structure within a factored representation.

The online approach—using look-ahead search to select an action for the current belief state—was first examined by Satia and Lave (1973). The use of sampling at chance nodes was explored analytically by Kearns *et al.* (2000) and Ng and Jordan (2000). The basic ideas for an agent architecture using dynamic decision networks were proposed by Dean and Kanazawa (1989a). The book *Planning and Control* by Dean and Wellman (1991) goes

FACTORED MDP

RELATIONAL MDP

into much greater depth, making connections between DBN/DDN models and the classical control literature on filtering. Tatman and Shachter (1990) showed how to apply dynamic programming algorithms to DDN models. Russell (1998) explains various ways in which such agents can be scaled up and identifies a number of open research issues.

The roots of game theory can be traced back to proposals made in the 17th century by Christiaan Huygens and Gottfried Leibniz to study competitive and cooperative human interactions scientifically and mathematically. Throughout the 19th century, several leading economists created simple mathematical examples to analyze particular examples of competitive situations. The first formal results in game theory are due to Zermelo (1913) (who had, the year before, suggested a form of minimax search for games, albeit an incorrect one). Emile Borel (1921) introduced the notion of a mixed strategy. John von Neumann (1928) proved that every two-person, zero-sum game has a maximin equilibrium in mixed strategies and a well-defined value. Von Neumann's collaboration with the economist Oskar Morgenstern led to the publication in 1944 of the *Theory of Games and Economic Behavior*, the defining book for game theory. Publication of the book was delayed by the wartime paper shortage until a member of the Rockefeller family personally subsidized its publication.

In 1950, at the age of 21, John Nash published his ideas concerning equilibria in general (non-zero-sum) games. His definition of an equilibrium solution, although originating in the work of Cournot (1838), became known as Nash equilibrium. After a long delay because of the schizophrenia he suffered from 1959 onward, Nash was awarded the Nobel Memorial Prize in Economics (along with Reinhard Selten and John Harsanyi) in 1994. The Bayes–Nash equilibrium is described by Harsanyi (1967) and discussed by Kadane and Larkey (1982). Some issues in the use of game theory for agent control are covered by Binmore (1982).

The prisoner's dilemma was invented as a classroom exercise by Albert W. Tucker in 1950 (based on an example by Merrill Flood and Melvin Dresher) and is covered extensively by Axelrod (1985) and Poundstone (1993). Repeated games were introduced by Luce and Raiffa (1957), and games of partial information in extensive form by Kuhn (1953). The first practical algorithm for sequential, partial-information games was developed within AI by Koller *et al.* (1996); the paper by Koller and Pfeffer (1997) provides a readable introduction to the field and describe a working system for representing and solving sequential games.

The use of abstraction to reduce a game tree to a size that can be solved with Koller's technique is discussed by Billings *et al.* (2003). Bowling *et al.* (2008) show how to use importance sampling to get a better estimate of the value of a strategy. Waugh *et al.* (2009) show that the abstraction approach is vulnerable to making systematic errors in approximating the equilibrium solution, meaning that the whole approach is on shaky ground: it works for some games but not others. Korb *et al.* (1999) experiment with an opponent model in the form of a Bayesian network. It plays five-card stud about as well as experienced humans. (Zinkevich *et al.*, 2008) show how an approach that minimizes regret can find approximate equilibria for abstractions with  $10^{12}$  states, 100 times more than previous methods.

Game theory and MDPs are combined in the theory of Markov games, also called stochastic games (Littman, 1994; Hu and Wellman, 1998). Shapley (1953) actually described the value iteration algorithm independently of Bellman, but his results were not widely appreciated, perhaps because they were presented in the context of Markov games. Evolu-

tionary game theory (Smith, 1982; Weibull, 1995) looks at strategy drift over time: if your opponent's strategy is changing, how should you react? Textbooks on game theory from an economics point of view include those by Myerson (1991), Fudenberg and Tirole (1991), Osborne (2004), and Osborne and Rubinstein (1994); Mailath and Samuelson (2006) concentrate on repeated games. From an AI perspective we have Nisan *et al.* (2007), Leyton-Brown and Shoham (2008), and Shoham and Leyton-Brown (2009).

The 2007 Nobel Memorial Prize in Economics went to Hurwicz, Maskin, and Myerson “for having laid the foundations of mechanism design theory” (Hurwicz, 1973). The tragedy of the commons, a motivating problem for the field, was presented by Hardin (1968). The revelation principle is due to Myerson (1986), and the revenue equivalence theorem was developed independently by Myerson (1981) and Riley and Samuelson (1981). Two economists, Milgrom (1997) and Klemperer (2002), write about the multibillion-dollar spectrum auctions they were involved in.

Mechanism design is used in multiagent planning (Hunsberger and Grosz, 2000; Stone *et al.*, 2009) and scheduling (Rassenti *et al.*, 1982). Varian (1995) gives a brief overview with connections to the computer science literature, and Rosenschein and Zlotkin (1994) present a book-length treatment with applications to distributed AI. Related work on distributed AI also goes under other names, including collective intelligence (Tumer and Wolpert, 2000; Segaran, 2007) and market-based control (Clearwater, 1996). Since 2001 there has been an annual Trading Agents Competition (TAC), in which agents try to make the best profit on a series of auctions (Wellman *et al.*, 2001; Arunachalam and Sadeh, 2005). Papers on computational issues in auctions often appear in the ACM Conferences on Electronic Commerce.

---

## EXERCISES

**17.1** For the  $4 \times 3$  world shown in Figure 17.1, calculate which squares can be reached from (1,1) by the action sequence  $[Up, Up, Right, Right, Right]$  and with what probabilities. Explain how this computation is related to the prediction task (see Section 15.2.1) for a hidden Markov model.

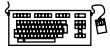
**17.2** Select a specific member of the set of policies that are optimal for  $R(s) > 0$  as shown in Figure 17.2(b), and calculate the fraction of time the agent spends in each state, in the limit, if the policy is executed forever. (*Hint:* Construct the state-to-state transition probability matrix corresponding to the policy and see Exercise 15.2.)

**17.3** Suppose that we define the utility of a state sequence to be the *maximum* reward obtained in any state in the sequence. Show that this utility function does not result in stationary preferences between state sequences. Is it still possible to define a utility function on states such that MEU decision making gives optimal behavior?

**17.4** Sometimes MDPs are formulated with a reward function  $R(s, a)$  that depends on the action taken or with a reward function  $R(s, a, s')$  that also depends on the outcome state.

- a. Write the Bellman equations for these formulations.

- b. Show how an MDP with reward function  $R(s, a, s')$  can be transformed into a different MDP with reward function  $R(s, a)$ , such that optimal policies in the new MDP correspond exactly to optimal policies in the original MDP.
- c. Now do the same to convert MDPs with  $R(s, a)$  into MDPs with  $R(s)$ .



**17.5** For the environment shown in Figure 17.1, find all the threshold values for  $R(s)$  such that the optimal policy changes when the threshold is crossed. You will need a way to calculate the optimal policy and its value for fixed  $R(s)$ . (*Hint*: Prove that the value of any fixed policy varies linearly with  $R(s)$ .)

**17.6** Equation (17.7) on page 654 states that the Bellman operator is a contraction.

- a. Show that, for any functions  $f$  and  $g$ ,

$$|\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|.$$

- b. Write out an expression for  $|(BU_i - BU'_i)(s)|$  and then apply the result from (a) to complete the proof that the Bellman operator is a contraction.

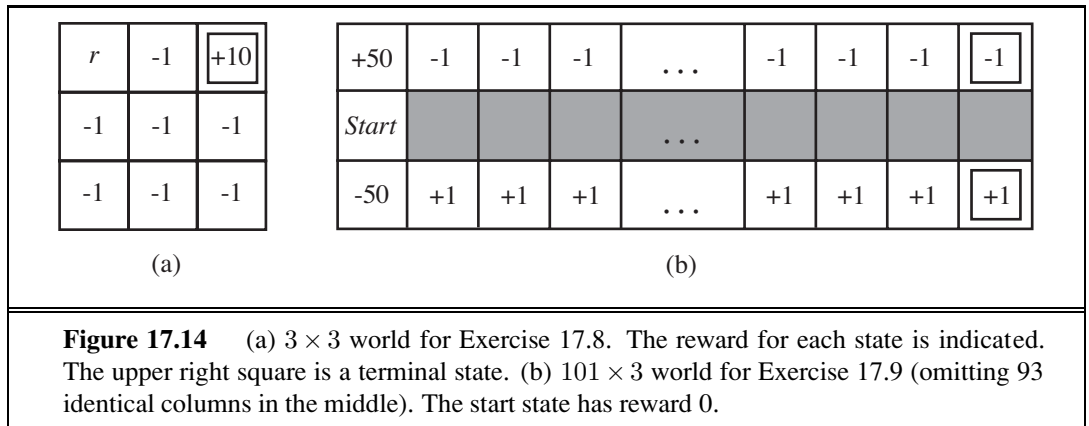
**17.7** This exercise considers two-player MDPs that correspond to zero-sum, turn-taking games like those in Chapter 5. Let the players be  $A$  and  $B$ , and let  $R(s)$  be the reward for player  $A$  in state  $s$ . (The reward for  $B$  is always equal and opposite.)

- a. Let  $U_A(s)$  be the utility of state  $s$  when it is  $A$ 's turn to move in  $s$ , and let  $U_B(s)$  be the utility of state  $s$  when it is  $B$ 's turn to move in  $s$ . All rewards and utilities are calculated from  $A$ 's point of view (just as in a minimax game tree). Write down Bellman equations defining  $U_A(s)$  and  $U_B(s)$ .
- b. Explain how to do two-player value iteration with these equations, and define a suitable termination criterion.
- c. Consider the game described in Figure 5.17 on page 197. Draw the state space (rather than the game tree), showing the moves by  $A$  as solid lines and moves by  $B$  as dashed lines. Mark each state with  $R(s)$ . You will find it helpful to arrange the states  $(s_A, s_B)$  on a two-dimensional grid, using  $s_A$  and  $s_B$  as “coordinates.”
- d. Now apply two-player value iteration to solve this game, and derive the optimal policy.

**17.8** Consider the  $3 \times 3$  world shown in Figure 17.14(a). The transition model is the same as in the  $4 \times 3$  Figure 17.1: 80% of the time the agent goes in the direction it selects; the rest of the time it moves at right angles to the intended direction.

Implement value iteration for this world for each value of  $r$  below. Use discounted rewards with a discount factor of 0.99. Show the policy obtained in each case. Explain intuitively why the value of  $r$  leads to each policy.

- a.  $r = 100$
- b.  $r = -3$
- c.  $r = 0$
- d.  $r = +3$



**17.9** Consider the  $101 \times 3$  world shown in Figure 17.14(b). In the start state the agent has a choice of two deterministic actions, *Up* or *Down*, but in the other states the agent has one deterministic action, *Right*. Assuming a discounted reward function, for what values of the discount  $\gamma$  should the agent choose *Up* and for which *Down*? Compute the utility of each action as a function of  $\gamma$ . (Note that this simple example actually reflects many real-world situations in which one must weigh the value of an immediate action versus the potential continual long-term consequences, such as choosing to dump pollutants into a lake.)

**17.10** Consider an undiscounted MDP having three states, (1, 2, 3), with rewards  $-1, -2, 0$ , respectively. State 3 is a terminal state. In states 1 and 2 there are two possible actions:  $a$  and  $b$ . The transition model is as follows:

- In state 1, action  $a$  moves the agent to state 2 with probability 0.8 and makes the agent stay put with probability 0.2.
- In state 2, action  $a$  moves the agent to state 1 with probability 0.8 and makes the agent stay put with probability 0.2.
- In either state 1 or state 2, action  $b$  moves the agent to state 3 with probability 0.1 and makes the agent stay put with probability 0.9.

Answer the following questions:

- a. What can be determined *qualitatively* about the optimal policy in states 1 and 2?
- b. Apply policy iteration, showing each step in full, to determine the optimal policy and the values of states 1 and 2. Assume that the initial policy has action  $b$  in both states.
- c. What happens to policy iteration if the initial policy has action  $a$  in both states? Does discounting help? Does the optimal policy depend on the discount factor?



**17.11** Consider the  $4 \times 3$  world shown in Figure 17.1.

- a. Implement an environment simulator for this environment, such that the specific geography of the environment is easily altered. Some code for doing this is already in the online code repository.

- b. Create an agent that uses policy iteration, and measure its performance in the environment simulator from various starting states. Perform several experiments from each starting state, and compare the average total reward received per run with the utility of the state, as determined by your algorithm.
- c. Experiment with increasing the size of the environment. How does the run time for policy iteration vary with the size of the environment?

**17.12** How can the value determination algorithm be used to calculate the expected loss experienced by an agent using a given set of utility estimates  $U$  and an estimated model  $P$ , compared with an agent using correct values?

**17.13** Let the initial belief state  $b_0$  for the  $4 \times 3$  POMDP on page 658 be the uniform distribution over the nonterminal states, i.e.,  $\langle \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, 0, 0 \rangle$ . Calculate the exact belief state  $b_1$  after the agent moves *Left* and its sensor reports 1 adjacent wall. Also calculate  $b_2$  assuming that the same thing happens again.

**17.14** What is the time complexity of  $d$  steps of POMDP value iteration for a sensorless environment?

**17.15** Consider a version of the two-state POMDP on page 661 in which the sensor is 90% reliable in state 0 but provides no information in state 1 (that is, it reports 0 or 1 with equal probability). Analyze, either qualitatively or quantitatively, the utility function and the optimal policy for this problem.

**17.16** Show that a dominant strategy equilibrium is a Nash equilibrium, but not vice versa.

**17.17** In the children's game of rock–paper–scissors each player reveals at the same time a choice of rock, paper, or scissors. Paper wraps rock, rock blunts scissors, and scissors cut paper. In the extended version rock–paper–scissors–fire–water, fire beats rock, paper, and scissors; rock, paper, and scissors beat water; and water beats fire. Write out the payoff matrix and find a mixed-strategy solution to this game.

**17.18** The following payoff matrix, from Blinder (1983) by way of Bernstein (1996), shows a game between politicians and the Federal Reserve.

	Fed: contract	Fed: do nothing	Fed: expand
Pol: contract	$F = 7, P = 1$	$F = 9, P = 4$	$F = 6, P = 6$
Pol: do nothing	$F = 8, P = 2$	$F = 5, P = 5$	$F = 4, P = 9$
Pol: expand	$F = 3, P = 3$	$F = 2, P = 7$	$F = 1, P = 8$

Politicians can expand or contract fiscal policy, while the Fed can expand or contract monetary policy. (And of course either side can choose to do nothing.) Each side also has preferences for who should do what—neither side wants to look like the bad guys. The payoffs shown are simply the rank orderings: 9 for first choice through 1 for last choice. Find the Nash equilibrium of the game in pure strategies. Is this a Pareto-optimal solution? You might wish to analyze the policies of recent administrations in this light.

**17.19** A Dutch auction is similar in an English auction, but rather than starting the bidding at a low price and increasing, in a Dutch auction the seller starts at a high price and gradually lowers the price until some buyer is willing to accept that price. (If multiple bidders accept the price, one is arbitrarily chosen as the winner.) More formally, the seller begins with a price  $p$  and gradually lowers  $p$  by increments of  $d$  until at least one buyer accepts the price. Assuming all bidders act rationally, is it true that for arbitrarily small  $d$ , a Dutch auction will always result in the bidder with the highest value for the item obtaining the item? If so, show mathematically why. If not, explain how it may be possible for the bidder with highest value for the item not to obtain it.

**17.20** Imagine an auction mechanism that is just like an ascending-bid auction, except that at the end, the winning bidder, the one who bid  $b_{max}$ , pays only  $b_{max}/2$  rather than  $b_{max}$ . Assuming all agents are rational, what is the expected revenue to the auctioneer for this mechanism, compared with a standard ascending-bid auction?

**17.21** Teams in the National Hockey League historically received 2 points for winning a game and 0 for losing. If the game is tied, an overtime period is played; if nobody wins in overtime, the game is a tie and each team gets 1 point. But league officials felt that teams were playing too conservatively in overtime (to avoid a loss), and it would be more exciting if overtime produced a winner. So in 1999 the officials experimented in mechanism design: the rules were changed, giving a team that loses in overtime 1 point, not 0. It is still 2 points for a win and 1 for a tie.

- a. Was hockey a zero-sum game before the rule change? After?
- b. Suppose that at a certain time  $t$  in a game, the home team has probability  $p$  of winning in regulation time, probability  $0.78 - p$  of losing, and probability 0.22 of going into overtime, where they have probability  $q$  of winning,  $.9 - q$  of losing, and .1 of tying. Give equations for the expected value for the home and visiting teams.
- c. Imagine that it were legal and ethical for the two teams to enter into a pact where they agree that they will skate to a tie in regulation time, and then both try in earnest to win in overtime. Under what conditions, in terms of  $p$  and  $q$ , would it be rational for both teams to agree to this pact?
- d. Longley and Sankaran (2005) report that since the rule change, the percentage of games with a winner in overtime went up 18.2%, as desired, but the percentage of overtime games also went up 3.6%. What does that suggest about possible collusion or conservative play after the rule change?



# 18 LEARNING FROM EXAMPLES

*In which we describe agents that can improve their behavior through diligent study of their own experiences.*

## LEARNING

An agent is **learning** if it improves its performance on future tasks after making observations about the world. Learning can range from the trivial, as exhibited by jotting down a phone number, to the profound, as exhibited by Albert Einstein, who inferred a new theory of the universe. In this chapter we will concentrate on one class of learning problem, which seems restricted but actually has vast applicability: from a collection of input–output pairs, learn a function that predicts the output for new inputs.

Why would we want an agent to learn? If the design of the agent can be improved, why wouldn't the designers just program in that improvement to begin with? There are three main reasons. First, the designers cannot anticipate all possible situations that the agent might find itself in. For example, a robot designed to navigate mazes must learn the layout of each new maze it encounters. Second, the designers cannot anticipate all changes over time; a program designed to predict tomorrow's stock market prices must learn to adapt when conditions change from boom to bust. Third, sometimes human programmers have no idea how to program a solution themselves. For example, most people are good at recognizing the faces of family members, but even the best programmers are unable to program a computer to accomplish that task, except by using learning algorithms. This chapter first gives an overview of the various forms of learning, then describes one popular approach, decision-tree learning, in Section 18.3, followed by a theoretical analysis of learning in Sections 18.4 and 18.5. We look at various learning systems used in practice: linear models, nonlinear models (in particular, neural networks), nonparametric models, and support vector machines. Finally we show how ensembles of models can outperform a single model.

## 18.1 FORMS OF LEARNING

---

Any component of an agent can be improved by learning from data. The improvements, and the techniques used to make them, depend on four major factors:

- Which *component* is to be improved.

- What *prior knowledge* the agent already has.
- What *representation* is used for the data and the component.
- What *feedback* is available to learn from.

### Components to be learned

Chapter 2 described several agent designs. The components of these agents include:

1. A direct mapping from conditions on the current state to actions.
2. A means to infer relevant properties of the world from the percept sequence.
3. Information about the way the world evolves and about the results of possible actions the agent can take.
4. Utility information indicating the desirability of world states.
5. Action-value information indicating the desirability of actions.
6. Goals that describe classes of states whose achievement maximizes the agent's utility.

Each of these components can be learned. Consider, for example, an agent training to become a taxi driver. Every time the instructor shouts “Brake!” the agent might learn a condition–action rule for when to brake (component 1); the agent also learns every time the instructor does not shout. By seeing many camera images that it is told contain buses, it can learn to recognize them (2). By trying actions and observing the results—for example, braking hard on a wet road—it can learn the effects of its actions (3). Then, when it receives no tip from passengers who have been thoroughly shaken up during the trip, it can learn a useful component of its overall utility function (4).

### Representation and prior knowledge

We have seen several examples of representations for agent components: propositional and first-order logical sentences for the components in a logical agent; Bayesian networks for the inferential components of a decision-theoretic agent, and so on. Effective learning algorithms have been devised for all of these representations. This chapter (and most of current machine learning research) covers inputs that form a **factored representation**—a vector of attribute values—and outputs that can be either a continuous numerical value or a discrete value. Chapter 19 covers functions and prior knowledge composed of first-order logic sentences, and Chapter 20 concentrates on Bayesian networks.

There is another way to look at the various types of learning. We say that learning a (possibly incorrect) general function or rule from specific input–output pairs is called **inductive learning**. We will see in Chapter 19 that we can also do **analytical** or **deductive learning**: going from a known general rule to a new rule that is logically entailed, but is useful because it allows more efficient processing.

### Feedback to learn from

There are three *types of feedback* that determine the three main types of learning:

In **unsupervised learning** the agent learns patterns in the input even though no explicit feedback is supplied. The most common unsupervised learning task is **clustering**: detecting

potentially useful clusters of input examples. For example, a taxi agent might gradually develop a concept of “good traffic days” and “bad traffic days” without ever being given labeled examples of each by a teacher.

REINFORCEMENT  
LEARNING

In **reinforcement learning** the agent learns from a series of reinforcements—rewards or punishments. For example, the lack of a tip at the end of the journey gives the taxi agent an indication that it did something wrong. The two points for a win at the end of a chess game tells the agent it did something right. It is up to the agent to decide which of the actions prior to the reinforcement were most responsible for it.

SUPERVISED  
LEARNING

In **supervised learning** the agent observes some example input–output pairs and learns a function that maps from input to output. In component 1 above, the inputs are percepts and the output are provided by a teacher who says “Brake!” or “Turn left.” In component 2, the inputs are camera images and the outputs again come from a teacher who says “that’s a bus.” In 3, the theory of braking is a function from states and braking actions to stopping distance in feet. In this case the output value is available directly from the agent’s percepts (after the fact); the environment is the teacher.

SEMI-SUPERVISED  
LEARNING

In practice, these distinction are not always so crisp. In **semi-supervised learning** we are given a few labeled examples and must make what we can of a large collection of unlabeled examples. Even the labels themselves may not be the oracular truths that we hope for. Imagine that you are trying to build a system to guess a person’s age from a photo. You gather some labeled examples by snapping pictures of people and asking their age. That’s supervised learning. But in reality some of the people lied about their age. It’s not just that there is random noise in the data; rather the inaccuracies are systematic, and to uncover them is an unsupervised learning problem involving images, self-reported ages, and true (unknown) ages. Thus, both noise and lack of labels create a continuum between supervised and unsupervised learning.

## 18.2 SUPERVISED LEARNING

The task of supervised learning is this:

TRAINING SET

Given a **training set** of  $N$  example input–output pairs

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N),$$

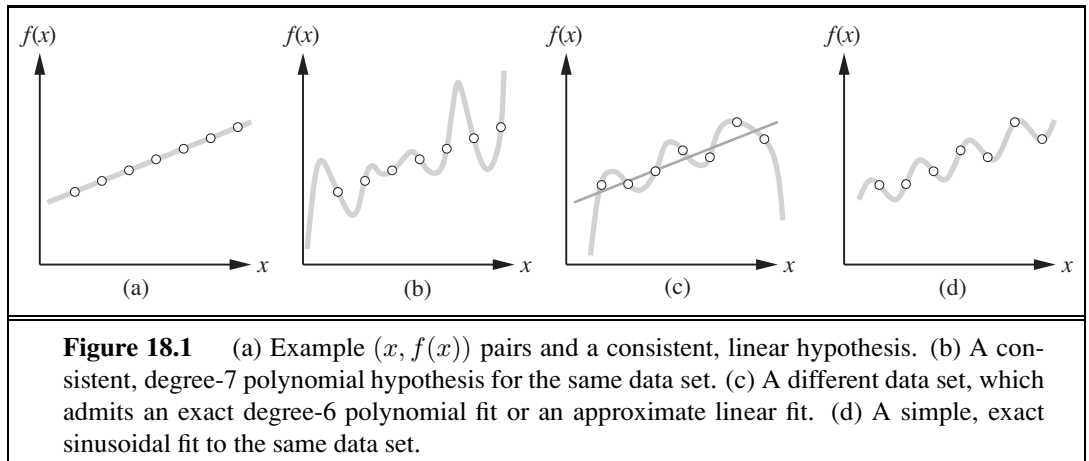
where each  $y_j$  was generated by an unknown function  $y = f(x)$ ,  
discover a function  $h$  that approximates the true function  $f$ .

HYPOTHESIS

Here  $x$  and  $y$  can be any value; they need not be numbers. The function  $h$  is a **hypothesis**.<sup>1</sup> Learning is a search through the space of possible hypotheses for one that will perform well, even on new examples beyond the training set. To measure the accuracy of a hypothesis we give it a **test set** of examples that are distinct from the training set. We say a hypothesis

TEST SET

<sup>1</sup> A note on notation: except where noted, we will use  $j$  to index the  $N$  examples;  $x_j$  will always be the input and  $y_j$  the output. In cases where the input is specifically a vector of attribute values (beginning with Section 18.3), we will use  $\mathbf{x}_j$  for the  $j$ th example and we will use  $i$  to index the  $n$  attributes of each example. The elements of  $\mathbf{x}_j$  are written  $x_{j,1}, x_{j,2}, \dots, x_{j,n}$ .



## GENERALIZATION

**generalizes** well if it correctly predicts the value of  $y$  for novel examples. Sometimes the function  $f$  is stochastic—it is not strictly a function of  $x$ , and what we have to learn is a conditional probability distribution,  $\mathbf{P}(Y | x)$ .

## CLASSIFICATION

When the output  $y$  is one of a finite set of values (such as *sunny*, *cloudy* or *rainy*), the learning problem is called **classification**, and is called Boolean or binary classification if there are only two values. When  $y$  is a number (such as tomorrow's temperature), the learning problem is called **regression**. (Technically, solving a regression problem is finding a conditional expectation or average value of  $y$ , because the probability that we have found *exactly* the right real-valued number for  $y$  is 0.)

## REGRESSION

## HYPOTHESIS SPACE

Figure 18.1 shows a familiar example: fitting a function of a single variable to some data points. The examples are points in the  $(x, y)$  plane, where  $y = f(x)$ . We don't know what  $f$  is, but we will approximate it with a function  $h$  selected from a **hypothesis space**,  $\mathcal{H}$ , which for this example we will take to be the set of polynomials, such as  $x^5 + 3x^2 + 2$ . Figure 18.1(a) shows some data with an exact fit by a straight line (the polynomial  $0.4x + 3$ ). The line is called a **consistent** hypothesis because it agrees with all the data. Figure 18.1(b) shows a high-degree polynomial that is also consistent with the same data. This illustrates a fundamental problem in inductive learning: *how do we choose from among multiple consistent hypotheses?* One answer is to prefer the *simplest* hypothesis consistent with the data. This principle is called **Ockham's razor**, after the 14th-century English philosopher William of Ockham, who used it to argue sharply against all sorts of complications. Defining simplicity is not easy, but it seems clear that a degree-1 polynomial is simpler than a degree-7 polynomial, and thus (a) should be preferred to (b). We will make this intuition more precise in Section 18.4.3.

## CONSISTENT



## OCKHAM'S RAZOR



Figure 18.1(c) shows a second data set. There is no consistent straight line for this data set; in fact, it requires a degree-6 polynomial for an exact fit. There are just 7 data points, so a polynomial with 7 parameters does not seem to be finding any pattern in the data and we do not expect it to generalize well. A straight line that is not consistent with any of the data points, but might generalize fairly well for unseen values of  $x$ , is also shown in (c). *In general, there is a tradeoff between complex hypotheses that fit the training data well and simpler hypotheses that may generalize better.* In Figure 18.1(d) we expand the

REALIZABLE

hypothesis space  $\mathcal{H}$  to allow polynomials over both  $x$  and  $\sin(x)$ , and find that the data in (c) can be fitted exactly by a simple function of the form  $ax + b + c \sin(x)$ . This shows the importance of the choice of hypothesis space. We say that a learning problem is **realizable** if the hypothesis space contains the true function. Unfortunately, we cannot always tell whether a given learning problem is realizable, because the true function is not known.

In some cases, an analyst looking at a problem is willing to make more fine-grained distinctions about the hypothesis space, to say—even before seeing any data—not just that a hypothesis is possible or impossible, but rather how probable it is. Supervised learning can be done by choosing the hypothesis  $h^*$  that is most probable given the data:

$$h^* = \operatorname{argmax}_{h \in \mathcal{H}} P(h|data) .$$

By Bayes' rule this is equivalent to

$$h^* = \operatorname{argmax}_{h \in \mathcal{H}} P(data|h) P(h) .$$

Then we can say that the prior probability  $P(h)$  is high for a degree-1 or -2 polynomial, lower for a degree-7 polynomial, and especially low for degree-7 polynomials with large, sharp spikes as in Figure 18.1(b). We allow unusual-looking functions when the data say we really need them, but we discourage them by giving them a low prior probability.



Why not let  $\mathcal{H}$  be the class of all Java programs, or Turing machines? After all, every computable function can be represented by some Turing machine, and that is the best we can do. One problem with this idea is that it does not take into account the computational complexity of learning. *There is a tradeoff between the expressiveness of a hypothesis space and the complexity of finding a good hypothesis within that space.* For example, fitting a straight line to data is an easy computation; fitting high-degree polynomials is somewhat harder; and fitting Turing machines is in general undecidable. A second reason to prefer simple hypothesis spaces is that presumably we will want to use  $h$  after we have learned it, and computing  $h(x)$  when  $h$  is a linear function is guaranteed to be fast, while computing an arbitrary Turing machine program is not even guaranteed to terminate. For these reasons, most work on learning has focused on simple representations.

We will see that the expressiveness–complexity tradeoff is not as simple as it first seems: it is often the case, as we saw with first-order logic in Chapter 8, that an expressive language makes it possible for a *simple* hypothesis to fit the data, whereas restricting the expressiveness of the language means that any consistent hypothesis must be very complex. For example, the rules of chess can be written in a page or two of first-order logic, but require thousands of pages when written in propositional logic.

## 18.3 LEARNING DECISION TREES

Decision tree induction is one of the simplest and yet most successful forms of machine learning. We first describe the representation—the hypothesis space—and then show how to learn a good hypothesis.

### 18.3.1 The decision tree representation

DECISION TREE

A **decision tree** represents a function that takes as input a vector of attribute values and returns a “decision”—a single output value. The input and output values can be discrete or continuous. For now we will concentrate on problems where the inputs have discrete values and the output has exactly two possible values; this is Boolean classification, where each example input will be classified as true (a **positive** example) or false (a **negative** example).

POSITIVE

NEGATIVE

A decision tree reaches its decision by performing a sequence of tests. Each internal node in the tree corresponds to a test of the value of one of the input attributes,  $A_i$ , and the branches from the node are labeled with the possible values of the attribute,  $A_i = v_{ik}$ . Each leaf node in the tree specifies a value to be returned by the function. The decision tree representation is natural for humans; indeed, many “How To” manuals (e.g., for car repair) are written entirely as a single decision tree stretching over hundreds of pages.

GOAL PREDICATE

As an example, we will build a decision tree to decide whether to wait for a table at a restaurant. The aim here is to learn a definition for the **goal predicate** *WillWait*. First we list the attributes that we will consider as part of the input:

1. *Alternate*: whether there is a suitable alternative restaurant nearby.
2. *Bar*: whether the restaurant has a comfortable bar area to wait in.
3. *Fri/Sat*: true on Fridays and Saturdays.
4. *Hungry*: whether we are hungry.
5. *Patrons*: how many people are in the restaurant (values are *None*, *Some*, and *Full*).
6. *Price*: the restaurant’s price range (\$, \$\$, \$\$\$).
7. *Raining*: whether it is raining outside.
8. *Reservation*: whether we made a reservation.
9. *Type*: the kind of restaurant (French, Italian, Thai, or burger).
10. *WaitEstimate*: the wait estimated by the host (0–10 minutes, 10–30, 30–60, or >60).

Note that every variable has a small set of possible values; the value of *WaitEstimate*, for example, is not an integer, rather it is one of the four discrete values 0–10, 10–30, 30–60, or >60. The decision tree usually used by one of us (SR) for this domain is shown in Figure 18.2. Notice that the tree ignores the *Price* and *Type* attributes. Examples are processed by the tree starting at the root and following the appropriate branch until a leaf is reached. For instance, an example with *Patrons* = *Full* and *WaitEstimate* = 0–10 will be classified as positive (i.e., yes, we will wait for a table).

### 18.3.2 Expressiveness of decision trees

A Boolean decision tree is logically equivalent to the assertion that the goal attribute is true if and only if the input attributes satisfy one of the paths leading to a leaf with value *true*. Writing this out in propositional logic, we have

$$Goal \Leftrightarrow (Path_1 \vee Path_2 \vee \dots),$$

where each *Path* is a conjunction of attribute-value tests required to follow that path. Thus, the whole expression is equivalent to disjunctive normal form (see page 283), which means

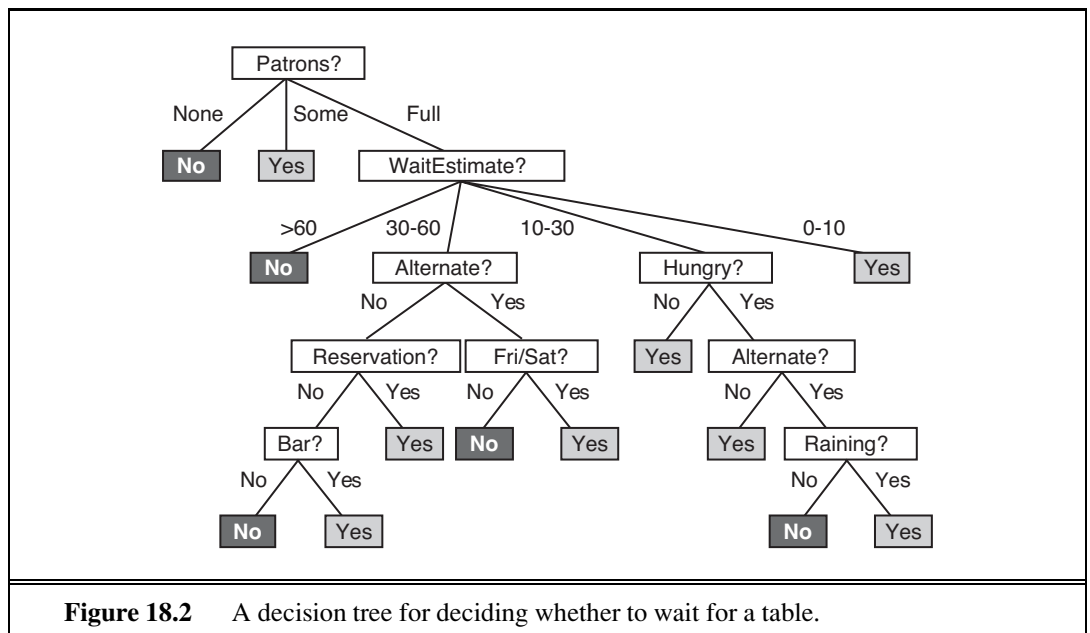
that any function in propositional logic can be expressed as a decision tree. As an example, the rightmost path in Figure 18.2 is

$$\text{Path} = (\text{Patrons} = \text{Full} \wedge \text{WaitEstimate} = 0-10) .$$

For a wide variety of problems, the decision tree format yields a nice, concise result. But some functions cannot be represented concisely. For example, the majority function, which returns true if and only if more than half of the inputs are true, requires an exponentially large decision tree. In other words, decision trees are good for some kinds of functions and bad for others. Is there *any* kind of representation that is efficient for *all* kinds of functions? Unfortunately, the answer is no. We can show this in a general way. Consider the set of all Boolean functions on  $n$  attributes. How many different functions are in this set? This is just the number of different truth tables that we can write down, because the function is defined by its truth table. A truth table over  $n$  attributes has  $2^n$  rows, one for each combination of values of the attributes. We can consider the “answer” column of the table as a  $2^n$ -bit number that defines the function. That means there are  $2^{2^n}$  different functions (and there will be more than that number of trees, since more than one tree can compute the same function). This is a scary number. For example, with just the ten Boolean attributes of our restaurant problem there are  $2^{1024}$  or about  $10^{308}$  different functions to choose from, and for 20 attributes there are over  $10^{300,000}$ . We will need some ingenious algorithms to find good hypotheses in such a large space.

### 18.3.3 Inducing decision trees from examples

An example for a Boolean decision tree consists of an  $(\mathbf{x}, y)$  pair, where  $\mathbf{x}$  is a vector of values for the input attributes, and  $y$  is a single Boolean output value. A training set of 12 examples



Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$x_1$	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>0-10</i>	$y_1 = \text{Yes}$
$x_2$	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>30-60</i>	$y_2 = \text{No}$
$x_3$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Some</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>0-10</i>	$y_3 = \text{Yes}$
$x_4$	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Thai</i>	<i>10-30</i>	$y_4 = \text{Yes}$
$x_5$	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>&gt;60</i>	$y_5 = \text{No}$
$x_6$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Italian</i>	<i>0-10</i>	$y_6 = \text{Yes}$
$x_7$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>0-10</i>	$y_7 = \text{No}$
$x_8$	<i>No</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Thai</i>	<i>0-10</i>	$y_8 = \text{Yes}$
$x_9$	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>&gt;60</i>	$y_9 = \text{No}$
$x_{10}$	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>Italian</i>	<i>10-30</i>	$y_{10} = \text{No}$
$x_{11}$	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>0-10</i>	$y_{11} = \text{No}$
$x_{12}$	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>30-60</i>	$y_{12} = \text{Yes}$

**Figure 18.3** Examples for the restaurant domain.

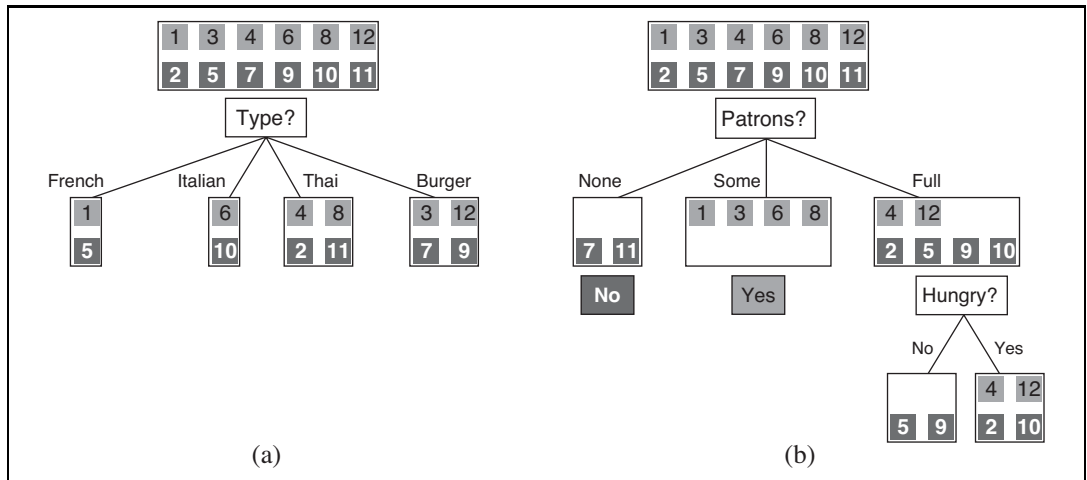
is shown in Figure 18.3. The positive examples are the ones in which the goal *WillWait* is true ( $x_1, x_3, \dots$ ); the negative examples are the ones in which it is false ( $x_2, x_5, \dots$ ).

We want a tree that is consistent with the examples and is as small as possible. Unfortunately, no matter how we measure size, it is an intractable problem to find the smallest consistent tree; there is no way to efficiently search through the  $2^{2^n}$  trees. With some simple heuristics, however, we can find a good approximate solution: a small (but not smallest) consistent tree. The DECISION-TREE-LEARNING algorithm adopts a greedy divide-and-conquer strategy: always test the most important attribute first. This test divides the problem up into smaller subproblems that can then be solved recursively. By “most important attribute,” we mean the one that makes the most difference to the classification of an example. That way, we hope to get to the correct classification with a small number of tests, meaning that all paths in the tree will be short and the tree as a whole will be shallow.

Figure 18.4(a) shows that *Type* is a poor attribute, because it leaves us with four possible outcomes, each of which has the same number of positive as negative examples. On the other hand, in (b) we see that *Patrons* is a fairly important attribute, because if the value is *None* or *Some*, then we are left with example sets for which we can answer definitively (*No* and *Yes*, respectively). If the value is *Full*, we are left with a mixed set of examples. In general, after the first attribute test splits up the examples, each outcome is a new decision tree learning problem in itself, with fewer examples and one less attribute. There are four cases to consider for these recursive problems:

1. If the remaining examples are all positive (or all negative), then we are done: we can answer *Yes* or *No*. Figure 18.4(b) shows examples of this happening in the *None* and *Some* branches.
2. If there are some positive and some negative examples, then choose the best attribute to split them. Figure 18.4(b) shows *Hungry* being used to split the remaining examples.
3. If there are no examples left, it means that no example has been observed for this com-





**Figure 18.4** Splitting the examples by testing on attributes. At each node we show the positive (light boxes) and negative (dark boxes) examples remaining. (a) Splitting on *Type* brings us no nearer to distinguishing between positive and negative examples. (b) Splitting on *Patrons* does a good job of separating positive and negative examples. After splitting on *Patrons*, *Hungry* is a fairly good second test.

bination of attribute values, and we return a default value calculated from the plurality classification of all the examples that were used in constructing the node's parent. These are passed along in the variable *parent\_examples*.

4. If there are no attributes left, but both positive and negative examples, it means that these examples have exactly the same description, but different classifications. This can happen because there is an error or **noise** in the data; because the domain is nondeterministic; or because we can't observe an attribute that would distinguish the examples. The best we can do is return the plurality classification of the remaining examples.

NOISE

The DECISION-TREE-LEARNING algorithm is shown in Figure 18.5. Note that the set of examples is crucial for *constructing* the tree, but nowhere do the examples appear in the tree itself. A tree consists of just tests on attributes in the interior nodes, values of attributes on the branches, and output values on the leaf nodes. The details of the IMPORTANCE function are given in Section 18.3.4. The output of the learning algorithm on our sample training set is shown in Figure 18.6. The tree is clearly different from the original tree shown in Figure 18.2. One might conclude that the learning algorithm is not doing a very good job of learning the correct function. This would be the wrong conclusion to draw, however. The learning algorithm looks at the *examples*, not at the correct function, and in fact, its hypothesis (see Figure 18.6) not only is consistent with all the examples, but is considerably simpler than the original tree! The learning algorithm has no reason to include tests for *Raining* and *Reservation*, because it can classify all the examples without them. It has also detected an interesting and previously unsuspected pattern: the first author will wait for Thai food on weekends. It is also bound to make some mistakes for cases where it has seen no examples. For example, it has never seen a case where the wait is 0–10 minutes but the restaurant is full.

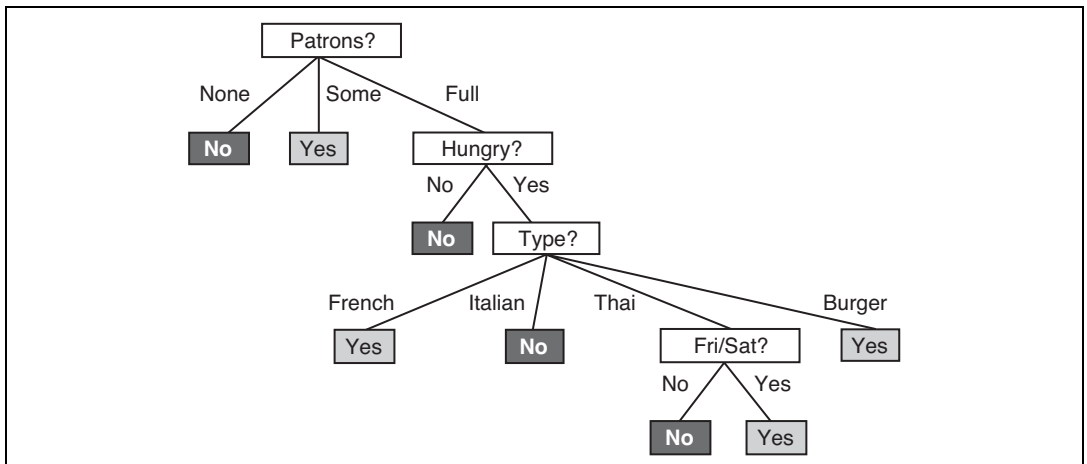
```

function DECISION-TREE-LEARNING(examples, attributes, parent_examples) returns
  a tree

  if examples is empty then return PLURALITY-VALUE(parent_examples)
  else if all examples have the same classification then return the classification
  else if attributes is empty then return PLURALITY-VALUE(examples)
  else
     $A \leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$ 
    tree  $\leftarrow$  a new decision tree with root test A
    for each value  $v_k$  of A do
      exs  $\leftarrow \{e : e \in \text{examples} \text{ and } e.A = v_k\}$ 
      subtree  $\leftarrow$  DECISION-TREE-LEARNING(exs, attributes - A, examples)
      add a branch to tree with label (A =  $v_k$ ) and subtree subtree
    return tree

```

**Figure 18.5** The decision-tree learning algorithm. The function IMPORTANCE is described in Section 18.3.4. The function PLURALITY-VALUE selects the most common output value among a set of examples, breaking ties randomly.

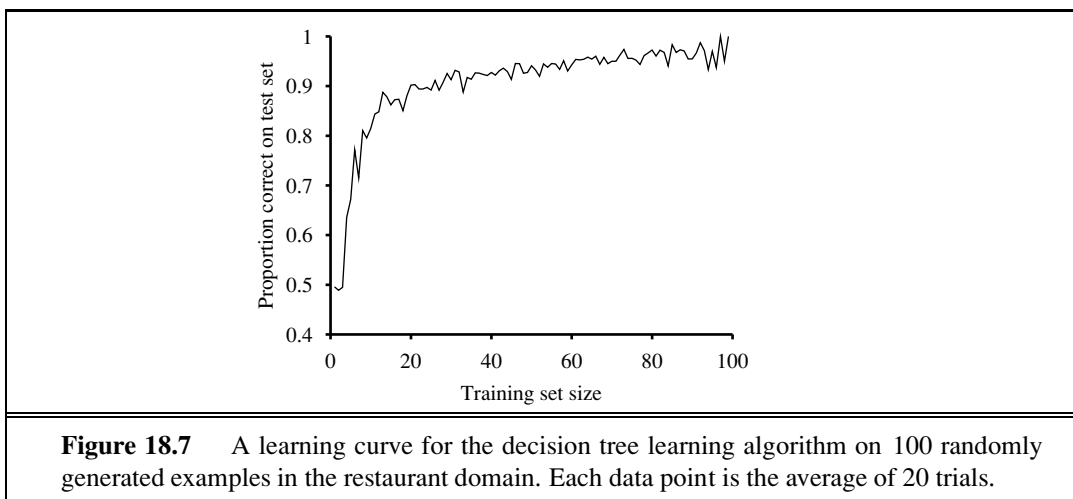


**Figure 18.6** The decision tree induced from the 12-example training set.

In that case it says not to wait when *Hungry* is false, but I (SR) would certainly wait. With more training examples the learning program could correct this mistake.

We note there is a danger of over-interpreting the tree that the algorithm selects. When there are several variables of similar importance, the choice between them is somewhat arbitrary: with slightly different input examples, a different variable would be chosen to split on first, and the whole tree would look completely different. The function computed by the tree would still be similar, but the structure of the tree can vary widely.

We can evaluate the accuracy of a learning algorithm with a **learning curve**, as shown in Figure 18.7. We have 100 examples at our disposal, which we split into a training set and



a test set. We learn a hypothesis  $h$  with the training set and measure its accuracy with the test set. We do this starting with a training set of size 1 and increasing one at a time up to size 99. For each size we actually repeat the process of randomly splitting 20 times, and average the results of the 20 trials. The curve shows that as the training set size grows, the accuracy increases. (For this reason, learning curves are also called **happy graphs**.) In this graph we reach 95% accuracy, and it looks like the curve might continue to increase with more data.

### 18.3.4 Choosing attribute tests

The greedy search used in decision tree learning is designed to approximately minimize the depth of the final tree. The idea is to pick the attribute that goes as far as possible toward providing an exact classification of the examples. A perfect attribute divides the examples into sets, each of which are all positive or all negative and thus will be leaves of the tree. The *Patrons* attribute is not perfect, but it is fairly good. A really useless attribute, such as *Type*, leaves the example sets with roughly the same proportion of positive and negative examples as the original set.

All we need, then, is a formal measure of “fairly good” and “really useless” and we can implement the IMPORTANCE function of Figure 18.5. We will use the notion of information gain, which is defined in terms of **entropy**, the fundamental quantity in information theory (Shannon and Weaver, 1949).

Entropy is a measure of the uncertainty of a random variable; acquisition of information corresponds to a reduction in entropy. A random variable with only one value—a coin that always comes up heads—has no uncertainty and thus its entropy is defined as zero; thus, we gain no information by observing its value. A flip of a fair coin is equally likely to come up heads or tails, 0 or 1, and we will soon show that this counts as “1 bit” of entropy. The roll of a fair *four*-sided die has 2 bits of entropy, because it takes two bits to describe one of four equally probable choices. Now consider an unfair coin that comes up heads 99% of the time. Intuitively, this coin has less uncertainty than the fair coin—if we guess heads we’ll be wrong only 1% of the time—so we would like it to have an entropy measure that is close to zero, but

positive. In general, the entropy of a random variable  $V$  with values  $v_k$ , each with probability  $P(v_k)$ , is defined as

$$\text{Entropy: } H(V) = \sum_k P(v_k) \log_2 \frac{1}{P(v_k)} = - \sum_k P(v_k) \log_2 P(v_k) .$$

We can check that the entropy of a fair coin flip is indeed 1 bit:

$$H(\text{Fair}) = -(0.5 \log_2 0.5 + 0.5 \log_2 0.5) = 1 .$$

If the coin is loaded to give 99% heads, we get

$$H(\text{Loaded}) = -(0.99 \log_2 0.99 + 0.01 \log_2 0.01) \approx 0.08 \text{ bits} .$$

It will help to define  $B(q)$  as the entropy of a Boolean random variable that is true with probability  $q$ :

$$B(q) = -(q \log_2 q + (1 - q) \log_2 (1 - q)) .$$

Thus,  $H(\text{Loaded}) = B(0.99) \approx 0.08$ . Now let's get back to decision tree learning. If a training set contains  $p$  positive examples and  $n$  negative examples, then the entropy of the goal attribute on the whole set is

$$H(\text{Goal}) = B\left(\frac{p}{p+n}\right) .$$

The restaurant training set in Figure 18.3 has  $p = n = 6$ , so the corresponding entropy is  $B(0.5)$  or exactly 1 bit. A test on a single attribute  $A$  might give us only part of this 1 bit. We can measure exactly how much by looking at the entropy remaining *after* the attribute test.

An attribute  $A$  with  $d$  distinct values divides the training set  $E$  into subsets  $E_1, \dots, E_d$ . Each subset  $E_k$  has  $p_k$  positive examples and  $n_k$  negative examples, so if we go along that branch, we will need an additional  $B(p_k/(p_k + n_k))$  bits of information to answer the question. A randomly chosen example from the training set has the  $k$ th value for the attribute with probability  $(p_k + n_k)/(p + n)$ , so the expected entropy remaining after testing attribute  $A$  is

$$\text{Remainder}(A) = \sum_{k=1}^d \frac{p_k + n_k}{p + n} B\left(\frac{p_k}{p_k + n_k}\right) .$$

INFORMATION GAIN

The **information gain** from the attribute test on  $A$  is the expected reduction in entropy:

$$\text{Gain}(A) = B\left(\frac{p}{p+n}\right) - \text{Remainder}(A) .$$

In fact  $\text{Gain}(A)$  is just what we need to implement the IMPORTANCE function. Returning to the attributes considered in Figure 18.4, we have

$$\text{Gain}(\text{Patrons}) = 1 - \left[ \frac{2}{12} B\left(\frac{0}{2}\right) + \frac{4}{12} B\left(\frac{4}{4}\right) + \frac{6}{12} B\left(\frac{2}{6}\right) \right] \approx 0.541 \text{ bits},$$

$$\text{Gain}(\text{Type}) = 1 - \left[ \frac{2}{12} B\left(\frac{1}{2}\right) + \frac{2}{12} B\left(\frac{1}{2}\right) + \frac{4}{12} B\left(\frac{2}{4}\right) + \frac{4}{12} B\left(\frac{2}{4}\right) \right] = 0 \text{ bits},$$

confirming our intuition that *Patrons* is a better attribute to split on. In fact, *Patrons* has the maximum gain of any of the attributes and would be chosen by the decision-tree learning algorithm as the root.

### 18.3.5 Generalization and overfitting

On some problems, the DECISION-TREE-LEARNING algorithm will generate a large tree when there is actually no pattern to be found. Consider the problem of trying to predict whether the roll of a die will come up as 6 or not. Suppose that experiments are carried out with various dice and that the attributes describing each training example include the color of the die, its weight, the time when the roll was done, and whether the experimenters had their fingers crossed. If the dice are fair, the right thing to learn is a tree with a single node that says “no.” But the DECISION-TREE-LEARNING algorithm will seize on any pattern it can find in the input. If it turns out that there are 2 rolls of a 7-gram blue die with fingers crossed and they both come out 6, then the algorithm may construct a path that predicts 6 in that case. This problem is called **overfitting**. A general phenomenon, overfitting occurs with all types of learners, even when the target function is not at all random. In Figure 18.1(b) and (c), we saw polynomial functions overfitting the data. Overfitting becomes more likely as the hypothesis space and the number of input attributes grows, and less likely as we increase the number of training examples.

OVERFITTING

DECISION TREE  
PRUNING

For decision trees, a technique called **decision tree pruning** combats overfitting. Pruning works by eliminating nodes that are not clearly relevant. We start with a full tree, as generated by DECISION-TREE-LEARNING. We then look at a test node that has only leaf nodes as descendants. If the test appears to be irrelevant—detecting only noise in the data—then we eliminate the test, replacing it with a leaf node. We repeat this process, considering each test with only leaf descendants, until each one has either been pruned or accepted as is.

The question is, how do we detect that a node is testing an irrelevant attribute? Suppose we are at a node consisting of  $p$  positive and  $n$  negative examples. If the attribute is irrelevant, we would expect that it would split the examples into subsets that each have roughly the same proportion of positive examples as the whole set,  $p/(p+n)$ , and so the information gain will be close to zero.<sup>2</sup> Thus, the information gain is a good clue to irrelevance. Now the question is, how large a gain should we require in order to split on a particular attribute?

SIGNIFICANCE TEST

NULL HYPOTHESIS

We can answer this question by using a statistical **significance test**. Such a test begins by assuming that there is no underlying pattern (the so-called **null hypothesis**). Then the actual data are analyzed to calculate the extent to which they deviate from a perfect absence of pattern. If the degree of deviation is statistically unlikely (usually taken to mean a 5% probability or less), then that is considered to be good evidence for the presence of a significant pattern in the data. The probabilities are calculated from standard distributions of the amount of deviation one would expect to see in random sampling.

In this case, the null hypothesis is that the attribute is irrelevant and, hence, that the information gain for an infinitely large sample would be zero. We need to calculate the probability that, under the null hypothesis, a sample of size  $v = n + p$  would exhibit the observed deviation from the expected distribution of positive and negative examples. We can measure the deviation by comparing the actual numbers of positive and negative examples in

<sup>2</sup> The gain will be strictly positive except for the unlikely case where all the proportions are *exactly* the same. (See Exercise 18.5.)

each subset,  $p_k$  and  $n_k$ , with the expected numbers,  $\hat{p}_k$  and  $\hat{n}_k$ , assuming true irrelevance:

$$\hat{p}_k = p \times \frac{p_k + n_k}{p + n} \quad \hat{n}_k = n \times \frac{p_k + n_k}{p + n} .$$

A convenient measure of the total deviation is given by

$$\Delta = \sum_{k=1}^d \frac{(p_k - \hat{p}_k)^2}{\hat{p}_k} + \frac{(n_k - \hat{n}_k)^2}{\hat{n}_k} .$$

Under the null hypothesis, the value of  $\Delta$  is distributed according to the  $\chi^2$  (chi-squared) distribution with  $v - 1$  degrees of freedom. We can use a  $\chi^2$  table or a standard statistical library routine to see if a particular  $\Delta$  value confirms or rejects the null hypothesis. For example, consider the restaurant type attribute, with four values and thus three degrees of freedom. A value of  $\Delta = 7.82$  or more would reject the null hypothesis at the 5% level (and a value of  $\Delta = 11.35$  or more would reject at the 1% level). Exercise 18.8 asks you to extend the DECISION-TREE-LEARNING algorithm to implement this form of pruning, which is known as  $\chi^2$  **pruning**.

$\chi^2$  PRUNING

With pruning, noise in the examples can be tolerated. Errors in the example's label (e.g., an example  $(\mathbf{x}, \text{Yes})$  that should be  $(\mathbf{x}, \text{No})$ ) give a linear increase in prediction error, whereas errors in the descriptions of examples (e.g.,  $\text{Price} = \$$  when it was actually  $\text{Price} = \$\$$ ) have an asymptotic effect that gets worse as the tree shrinks down to smaller sets. Pruned trees perform significantly better than unpruned trees when the data contain a large amount of noise. Also, the pruned trees are often much smaller and hence easier to understand.

EARLY STOPPING

One final warning: You might think that  $\chi^2$  pruning and information gain look similar, so why not combine them using an approach called **early stopping**—have the decision tree algorithm stop generating nodes when there is no good attribute to split on, rather than going to all the trouble of generating nodes and then pruning them away. The problem with early stopping is that it stops us from recognizing situations where there is no one good attribute, but there are combinations of attributes that are informative. For example, consider the XOR function of two binary attributes. If there are roughly equal number of examples for all four combinations of input values, then neither attribute will be informative, yet the correct thing to do is to split on one of the attributes (it doesn't matter which one), and then at the second level we will get splits that are informative. Early stopping would miss this, but generate-and-then-prune handles it correctly.

### 18.3.6 Broadening the applicability of decision trees

In order to extend decision tree induction to a wider variety of problems, a number of issues must be addressed. We will briefly mention several, suggesting that a full understanding is best obtained by doing the associated exercises:

- **Missing data:** In many domains, not all the attribute values will be known for every example. The values might have gone unrecorded, or they might be too expensive to obtain. This gives rise to two problems: First, given a complete decision tree, how should one classify an example that is missing one of the test attributes? Second, how

should one modify the information-gain formula when some examples have unknown values for the attribute? These questions are addressed in Exercise 18.9.

GAIN RATIO

- **Multivalued attributes:** When an attribute has many possible values, the information gain measure gives an inappropriate indication of the attribute's usefulness. In the extreme case, an attribute such as *ExactTime* has a different value for every example, which means each subset of examples is a singleton with a unique classification, and the information gain measure would have its highest value for this attribute. But choosing this split first is unlikely to yield the best tree. One solution is to use the **gain ratio** (Exercise 18.10). Another possibility is to allow a Boolean test of the form  $A = v_k$ , that is, picking out just one of the possible values for an attribute, leaving the remaining values to possibly be tested later in the tree.

SPLIT POINT

- **Continuous and integer-valued input attributes:** Continuous or integer-valued attributes such as *Height* and *Weight*, have an infinite set of possible values. Rather than generate infinitely many branches, decision-tree learning algorithms typically find the **split point** that gives the highest information gain. For example, at a given node in the tree, it might be the case that testing on  $Weight > 160$  gives the most information. Efficient methods exist for finding good split points: start by sorting the values of the attribute, and then consider only split points that are between two examples in sorted order that have different classifications, while keeping track of the running totals of positive and negative examples on each side of the split point. Splitting is the most expensive part of real-world decision tree learning applications.

REGRESSION TREE

- **Continuous-valued output attributes:** If we are trying to predict a numerical output value, such as the price of an apartment, then we need a **regression tree** rather than a classification tree. A regression tree has at each leaf a linear function of some subset of numerical attributes, rather than a single value. For example, the branch for two-bedroom apartments might end with a linear function of square footage, number of bathrooms, and average income for the neighborhood. The learning algorithm must decide when to stop splitting and begin applying linear regression (see Section 18.6) over the attributes.

A decision-tree learning system for real-world applications must be able to handle all of these problems. Handling continuous-valued variables is especially important, because both physical and financial processes provide numerical data. Several commercial packages have been built that meet these criteria, and they have been used to develop thousands of fielded systems. In many areas of industry and commerce, decision trees are usually the first method tried when a classification method is to be extracted from a data set. One important property of decision trees is that it is possible for a human to understand the reason for the output of the learning algorithm. (Indeed, this is a *legal requirement* for financial decisions that are subject to anti-discrimination laws.) This is a property not shared by some other representations, such as neural networks.

## 18.4 EVALUATING AND CHOOSING THE BEST HYPOTHESIS

STATIONARITY  
ASSUMPTION

We want to learn a hypothesis that fits the future data best. To make that precise we need to define “future data” and “best.” We make the **stationarity assumption**: that there is a probability distribution over examples that remains stationary over time. Each example data point (before we see it) is a random variable  $E_j$  whose observed value  $e_j = (x_j, y_j)$  is sampled from that distribution, and is independent of the previous examples:

$$\mathbf{P}(E_j | E_{j-1}, E_{j-2}, \dots) = \mathbf{P}(E_j),$$

and each example has an identical prior probability distribution:

$$\mathbf{P}(E_j) = \mathbf{P}(E_{j-1}) = \mathbf{P}(E_{j-2}) = \dots$$

I.I.D.

Examples that satisfy these assumptions are called *independent and identically distributed* or **i.i.d.**. An i.i.d. assumption connects the past to the future; without some such connection, all bets are off—the future could be anything. (We will see later that learning can still occur if there are *slow* changes in the distribution.)

ERROR RATE

The next step is to define “best fit.” We define the **error rate** of a hypothesis as the proportion of mistakes it makes—the proportion of times that  $h(x) \neq y$  for an  $(x, y)$  example. Now, just because a hypothesis  $h$  has a low error rate on the training set does not mean that it will generalize well. A professor knows that an exam will not accurately evaluate students if they have already seen the exam questions. Similarly, to get an accurate evaluation of a hypothesis, we need to test it on a set of examples it has not seen yet. The simplest approach is the one we have seen already: randomly split the available data into a training set from which the learning algorithm produces  $h$  and a test set on which the accuracy of  $h$  is evaluated. This method, sometimes called **holdout cross-validation**, has the disadvantage that it fails to use all the available data; if we use half the data for the test set, then we are only training on half the data, and we may get a poor hypothesis. On the other hand, if we reserve only 10% of the data for the test set, then we may, by statistical chance, get a poor estimate of the actual accuracy.

HOLDOUT  
CROSS-VALIDATION

K-FOLD  
CROSS-VALIDATION

We can squeeze more out of the data and still get an accurate estimate using a technique called  **$k$ -fold cross-validation**. The idea is that each example serves double duty—as training data and test data. First we split the data into  $k$  equal subsets. We then perform  $k$  rounds of learning; on each round  $1/k$  of the data is held out as a test set and the remaining examples are used as training data. The average test set score of the  $k$  rounds should then be a better estimate than a single score. Popular values for  $k$  are 5 and 10—enough to give an estimate that is statistically likely to be accurate, at a cost of 5 to 10 times longer computation time. The extreme is  $k = n$ , also known as **leave-one-out cross-validation** or **LOOCV**.

LEAVE-ONE-OUT  
CROSS-VALIDATION  
LOOCV

PEEKING

Despite the best efforts of statistical methodologists, users frequently invalidate their results by inadvertently **peeking** at the test data. Peeking can happen like this: A learning algorithm has various “knobs” that can be twiddled to tune its behavior—for example, various different criteria for choosing the next attribute in decision tree learning. The researcher generates hypotheses for various different settings of the knobs, measures their error rates on the test set, and reports the error rate of the best hypothesis. Alas, peeking has occurred! The



reason is that the hypothesis was selected *on the basis of its test set error rate*, so information about the test set has leaked into the learning algorithm.

Peeking is a consequence of using test-set performance to both *choose* a hypothesis and *evaluate* it. The way to avoid this is to *really* hold the test set out—lock it away until you are completely done with learning and simply wish to obtain an independent evaluation of the final hypothesis. (And then, if you don't like the results . . . you have to obtain, and lock away, a completely new test set if you want to go back and find a better hypothesis.) If the test set is locked away, but you still want to measure performance on unseen data as a way of selecting a good hypothesis, then divide the available data (without the test set) into a training set and a **validation set**. The next section shows how to use validation sets to find a good tradeoff between hypothesis complexity and goodness of fit.

VALIDATION SET

### 18.4.1 Model selection: Complexity versus goodness of fit

In Figure 18.1 (page 696) we showed that higher-degree polynomials can fit the training data better, but when the degree is too high they will overfit, and perform poorly on validation data.

Choosing the degree of the polynomial is an instance of the problem of **model selection**. You can think of the task of finding the best hypothesis as two tasks: model selection defines the hypothesis space and then **optimization** finds the best hypothesis within that space.

MODEL SELECTION

OPTIMIZATION

In this section we explain how to select among models that are parameterized by *size*. For example, with polynomials we have *size* = 1 for linear functions, *size* = 2 for quadratics, and so on. For decision trees, the size could be the number of nodes in the tree. In all cases we want to find the value of the *size* parameter that best balances underfitting and overfitting to give the best test set accuracy.

An algorithm to perform model selection and optimization is shown in Figure 18.8. It is a **wrapper** that takes a learning algorithm as an argument (DECISION-TREE-LEARNING, for example). The wrapper enumerates models according to a parameter, *size*. For each size, it uses cross validation on *Learner* to compute the average error rate on the training and test sets. We start with the smallest, simplest models (which probably underfit the data), and iterate, considering more complex models at each step, until the models start to overfit. In Figure 18.9 we see typical curves: the training set error decreases monotonically (although there may in general be slight random variation), while the validation set error decreases at first, and then increases when the model begins to overfit. The cross-validation procedure picks the value of *size* with the lowest validation set error; the bottom of the U-shaped curve. We then generate a hypothesis of that *size*, using all the data (without holding out any of it). Finally, of course, we should evaluate the returned hypothesis on a separate test set.

WRAPPER

This approach requires that the learning algorithm accept a parameter, *size*, and deliver a hypothesis of that size. As we said, for decision tree learning, the size can be the number of nodes. We can modify DECISION-TREE-LEARNER so that it takes the number of nodes as an input, builds the tree breadth-first rather than depth-first (but at each level it still chooses the highest gain attribute first), and stops when it reaches the desired number of nodes.

**function** CROSS-VALIDATION-WRAPPER(*Learner*, *k*, *examples*) **returns** a hypothesis

**local variables:** *errT*, an array, indexed by *size*, storing training-set error rates  
*errV*, an array, indexed by *size*, storing validation-set error rates

**for** *size* = 1 to  $\infty$  **do**

*errT*[*size*], *errV*[*size*]  $\leftarrow$  CROSS-VALIDATION(*Learner*, *size*, *k*, *examples*)

**if** *errT* has converged **then do**

*best\_size*  $\leftarrow$  the value of *size* with minimum *errV*[*size*]

**return** *Learner*(*best\_size*, *examples*)

**function** CROSS-VALIDATION(*Learner*, *size*, *k*, *examples*) **returns** two values:  
 average training set error rate, average validation set error rate

*fold\_errT*  $\leftarrow$  0; *fold\_errV*  $\leftarrow$  0

**for** *fold* = 1 to *k* **do**

*training\_set*, *validation\_set*  $\leftarrow$  PARTITION(*examples*, *fold*, *k*)

*h*  $\leftarrow$  *Learner*(*size*, *training\_set*)

*fold\_errT*  $\leftarrow$  *fold\_errT* + ERROR-RATE(*h*, *training\_set*)

*fold\_errV*  $\leftarrow$  *fold\_errV* + ERROR-RATE(*h*, *validation\_set*)

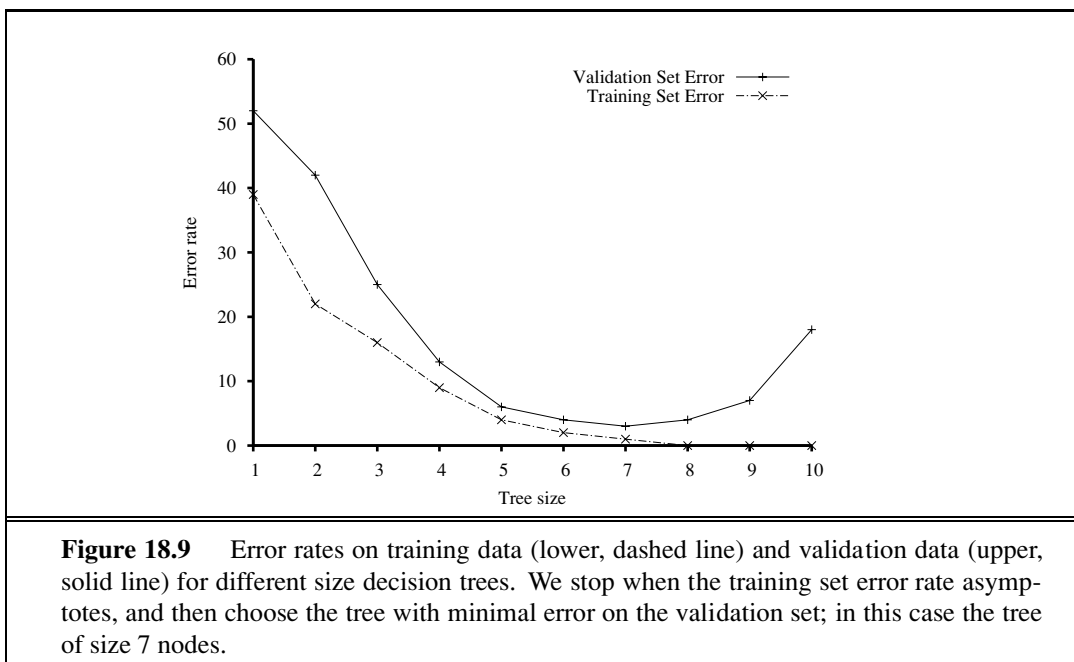
**return** *fold\_errT*/*k*, *fold\_errV*/*k*

**Figure 18.8** An algorithm to select the model that has the lowest error rate on validation data by building models of increasing complexity, and choosing the one with best empirical error rate on validation data. Here *errT* means error rate on the training data, and *errV* means error rate on the validation data. *Learner*(*size*, *examples*) returns a hypothesis whose complexity is set by the parameter *size*, and which is trained on the *examples*. PARTITION(*examples*, *fold*, *k*) splits *examples* into two subsets: a validation set of size  $N/k$  and a training set with all the other examples. The split is different for each value of *fold*.

### 18.4.2 From error rates to loss

So far, we have been trying to minimize error rate. This is clearly better than maximizing error rate, but it is not the full story. Consider the problem of classifying email messages as spam or non-spam. It is worse to classify non-spam as spam (and thus potentially miss an important message) than to classify spam as non-spam (and thus suffer a few seconds of annoyance). So a classifier with a 1% error rate, where almost all the errors were classifying spam as non-spam, would be better than a classifier with only a 0.5% error rate, if most of those errors were classifying non-spam as spam. We saw in Chapter 16 that decision-makers should maximize expected utility, and utility is what learners should maximize as well. In machine learning it is traditional to express utilities by means of a **loss function**. The loss function  $L(x, y, \hat{y})$  is defined as the amount of utility lost by predicting  $h(x) = \hat{y}$  when the correct answer is  $f(x) = y$ :

$$\begin{aligned} L(x, y, \hat{y}) &= \text{Utility}(\text{result of using } y \text{ given an input } x) \\ &\quad - \text{Utility}(\text{result of using } \hat{y} \text{ given an input } x) \end{aligned}$$



This is the most general formulation of the loss function. Often a simplified version is used,  $L(y, \hat{y})$ , that is independent of  $x$ . We will use the simplified version for the rest of this chapter, which means we can't say that it is worse to misclassify a letter from Mom than it is to misclassify a letter from our annoying cousin, but we can say it is 10 times worse to classify non-spam as spam than vice-versa:

$$L(\text{spam}, \text{nospam}) = 1, \quad L(\text{nospam}, \text{spam}) = 10.$$

Note that  $L(y, y)$  is always zero; by definition there is no loss when you guess exactly right. For functions with discrete outputs, we can enumerate a loss value for each possible misclassification, but we can't enumerate all the possibilities for real-valued data. If  $f(x)$  is 137.035999, we would be fairly happy with  $h(x) = 137.036$ , but just how happy should we be? In general small errors are better than large ones; two functions that implement that idea are the absolute value of the difference (called the  $L_1$  loss), and the square of the difference (called the  $L_2$  loss). If we are content with the idea of minimizing error rate, we can use the  $L_{0/1}$  loss function, which has a loss of 1 for an incorrect answer and is appropriate for discrete-valued outputs:

$$\text{Absolute value loss: } L_1(y, \hat{y}) = |y - \hat{y}|$$

$$\text{Squared error loss: } L_2(y, \hat{y}) = (y - \hat{y})^2$$

$$0/1 \text{ loss: } L_{0/1}(y, \hat{y}) = 0 \text{ if } y = \hat{y}, \text{ else } 1$$

The learning agent can theoretically maximize its expected utility by choosing the hypothesis that minimizes expected loss over all input-output pairs it will see. It is meaningless to talk about this expectation without defining a prior probability distribution,  $\mathbf{P}(X, Y)$  over examples. Let  $\mathcal{E}$  be the set of all possible input-output examples. Then the expected **generalization loss** for a hypothesis  $h$  (with respect to loss function  $L$ ) is

$$GenLoss_L(h) = \sum_{(x,y) \in \mathcal{E}} L(y, h(x)) P(x, y) ,$$

and the best hypothesis,  $h^*$ , is the one with the minimum expected generalization loss:

$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} GenLoss_L(h) .$$

Because  $P(x, y)$  is not known, the learning agent can only *estimate* generalization loss with **empirical loss** on a set of examples,  $E$ :

EMPIRICAL LOSS

$$EmpLoss_{L,E}(h) = \frac{1}{N} \sum_{(x,y) \in E} L(y, h(x)) .$$

The estimated best hypothesis  $\hat{h}^*$  is then the one with minimum empirical loss:

$$\hat{h}^* = \operatorname{argmin}_{h \in \mathcal{H}} EmpLoss_{L,E}(h) .$$

There are four reasons why  $\hat{h}^*$  may differ from the true function,  $f$ : unrealizability, variance, noise, and computational complexity. First,  $f$  may not be realizable—may not be in  $\mathcal{H}$ —or may be present in such a way that other hypotheses are preferred. Second, a learning algorithm will return different hypotheses for different sets of examples, even if those sets are drawn from the same true function  $f$ , and those hypotheses will make different predictions on new examples. The higher the variance among the predictions, the higher the probability of significant error. Note that even when the problem is realizable, there will still be random variance, but that variance decreases towards zero as the number of training examples increases. Third,  $f$  may be nondeterministic or **noisy**—it may return different values for  $f(x)$  each time  $x$  occurs. By definition, noise cannot be predicted; in many cases, it arises because the observed labels  $y$  are the result of attributes of the environment not listed in  $x$ . And finally, when  $\mathcal{H}$  is complex, it can be computationally intractable to systematically search the whole hypothesis space. The best we can do is a local search (hill climbing or greedy search) that explores only part of the space. That gives us an approximation error. Combining the sources of error, we're left with an estimation of an approximation of the true function  $f$ .

NOISE

Traditional methods in statistics and the early years of machine learning concentrated on **small-scale learning**, where the number of training examples ranged from dozens to the low thousands. Here the generalization error mostly comes from the approximation error of not having the true  $f$  in the hypothesis space, and from estimation error of not having enough training examples to limit variance. In recent years there has been more emphasis on **large-scale learning**, often with millions of examples. Here the generalization error is dominated by limits of computation: there is enough data and a rich enough model that we could find an  $h$  that is very close to the true  $f$ , but the computation to find it is too complex, so we settle for a sub-optimal approximation.

SMALL-SCALE  
LEARNINGLARGE-SCALE  
LEARNING

### 18.4.3 Regularization

In Section 18.4.1, we saw how to do model selection with cross-validation on model size. An alternative approach is to search for a hypothesis that directly minimizes the weighted sum of

empirical loss and the complexity of the hypothesis, which we will call the total cost:

$$\begin{aligned} \text{Cost}(h) &= \text{EmpLoss}(h) + \lambda \text{Complexity}(h) \\ \hat{h}^* &= \underset{h \in \mathcal{H}}{\operatorname{argmin}} \text{Cost}(h) . \end{aligned}$$

Here  $\lambda$  is a parameter, a positive number that serves as a conversion rate between loss and hypothesis complexity (which after all are not measured on the same scale). This approach combines loss and complexity into one metric, allowing us to find the best hypothesis all at once. Unfortunately we still need to do a cross-validation search to find the hypothesis that generalizes best, but this time it is with different values of  $\lambda$  rather than *size*. We select the value of  $\lambda$  that gives us the best validation set score.

REGULARIZATION

This process of explicitly penalizing complex hypotheses is called **regularization** (because it looks for a function that is more regular, or less complex). Note that the cost function requires us to make two choices: the loss function and the complexity measure, which is called a regularization function. The choice of regularization function depends on the hypothesis space. For example, a good regularization function for polynomials is the sum of the squares of the coefficients—keeping the sum small would guide us away from the wiggly polynomials in Figure 18.1(b) and (c). We will show an example of this type of regularization in Section 18.6.

FEATURE SELECTION

Another way to simplify models is to reduce the dimensions that the models work with. A process of **feature selection** can be performed to discard attributes that appear to be irrelevant.  $\chi^2$  pruning is a kind of feature selection.

MINIMUM  
DESCRIPTION  
LENGTH

It is in fact possible to have the empirical loss and the complexity measured on the same scale, without the conversion factor  $\lambda$ : they can both be measured in bits. First encode the hypothesis as a Turing machine program, and count the number of bits. Then count the number of bits required to encode the data, where a correctly predicted example costs zero bits and the cost of an incorrectly predicted example depends on how large the error is. The **minimum description length** or MDL hypothesis minimizes the total number of bits required. This works well in the limit, but for smaller problems there is a difficulty in that the choice of encoding for the program—for example, how best to encode a decision tree as a bit string—affects the outcome. In Chapter 20 (page 805), we describe a probabilistic interpretation of the MDL approach.

## 18.5 THE THEORY OF LEARNING

The main unanswered question in learning is this: How can we be sure that our learning algorithm has produced a hypothesis that will predict the correct value for previously unseen inputs? In formal terms, how do we know that the hypothesis  $h$  is close to the target function  $f$  if we don't know what  $f$  is? These questions have been pondered for several centuries. In more recent decades, other questions have emerged: how many examples do we need to get a good  $h$ ? What hypothesis space should we use? If the hypothesis space is very complex, can we even find the best  $h$ , or do we have to settle for a local maximum in the

space of hypotheses? How complex should  $h$  be? How do we avoid overfitting? This section examines these questions.

We'll start with the question of how many examples are needed for learning. We saw from the learning curve for decision tree learning on the restaurant problem (Figure 18.7 on page 703) that improves with more training data. Learning curves are useful, but they are specific to a particular learning algorithm on a particular problem. Are there some more general principles governing the number of examples needed in general? Questions like this are addressed by **computational learning theory**, which lies at the intersection of AI, statistics, and theoretical computer science. The underlying principle is that *any hypothesis that is seriously wrong will almost certainly be “found out” with high probability after a small number of examples, because it will make an incorrect prediction. Thus, any hypothesis that is consistent with a sufficiently large set of training examples is unlikely to be seriously wrong: that is, it must be **probably approximately correct***. Any learning algorithm that returns hypotheses that are probably approximately correct is called a **PAC learning** algorithm; we can use this approach to provide bounds on the performance of various learning algorithms.

PAC-learning theorems, like all theorems, are logical consequences of axioms. When a *theorem* (as opposed to, say, a political pundit) states something about the future based on the past, the axioms have to provide the “juice” to make that connection. For PAC learning, the juice is provided by the stationarity assumption introduced on page 708, which says that future examples are going to be drawn from the same fixed distribution  $\mathbf{P}(E) = \mathbf{P}(X, Y)$  as past examples. (Note that we do not have to know what distribution that is, just that it doesn't change.) In addition, to keep things simple, we will assume that the true function  $f$  is deterministic and is a member of the hypothesis class  $\mathcal{H}$  that is being considered.

The simplest PAC theorems deal with Boolean functions, for which the 0/1 loss is appropriate. The **error rate** of a hypothesis  $h$ , defined informally earlier, is defined formally here as the expected generalization error for examples drawn from the stationary distribution:

$$\text{error}(h) = \text{GenLoss}_{L_{0/1}}(h) = \sum_{x,y} L_{0/1}(y, h(x)) P(x, y) .$$

In other words,  $\text{error}(h)$  is the probability that  $h$  misclassifies a new example. This is the same quantity being measured experimentally by the learning curves shown earlier.

A hypothesis  $h$  is called **approximately correct** if  $\text{error}(h) \leq \epsilon$ , where  $\epsilon$  is a small constant. We will show that we can find an  $N$  such that, after seeing  $N$  examples, with high probability, all consistent hypotheses will be approximately correct. One can think of an approximately correct hypothesis as being “close” to the true function in hypothesis space: it lies inside what is called the  **$\epsilon$ -ball** around the true function  $f$ . The hypothesis space outside this ball is called  $\mathcal{H}_{\text{bad}}$ .

We can calculate the probability that a “seriously wrong” hypothesis  $h_b \in \mathcal{H}_{\text{bad}}$  is consistent with the first  $N$  examples as follows. We know that  $\text{error}(h_b) > \epsilon$ . Thus, the probability that it agrees with a given example is at most  $1 - \epsilon$ . Since the examples are independent, the bound for  $N$  examples is

$$P(h_b \text{ agrees with } N \text{ examples}) \leq (1 - \epsilon)^N .$$

COMPUTATIONAL  
LEARNING THEORY



PROBABLY  
APPROXIMATELY  
CORRECT  
PAC LEARNING

$\epsilon$ -BALL

The probability that  $\mathcal{H}_{\text{bad}}$  contains at least one consistent hypothesis is bounded by the sum of the individual probabilities:

$$P(\mathcal{H}_{\text{bad}} \text{ contains a consistent hypothesis}) \leq |\mathcal{H}_{\text{bad}}|(1 - \epsilon)^N \leq |\mathcal{H}|(1 - \epsilon)^N,$$

where we have used the fact that  $|\mathcal{H}_{\text{bad}}| \leq |\mathcal{H}|$ . We would like to reduce the probability of this event below some small number  $\delta$ :

$$|\mathcal{H}|(1 - \epsilon)^N \leq \delta.$$

Given that  $1 - \epsilon \leq e^{-\epsilon}$ , we can achieve this if we allow the algorithm to see

$$N \geq \frac{1}{\epsilon} \left( \ln \frac{1}{\delta} + \ln |\mathcal{H}| \right) \quad (18.1)$$

examples. Thus, if a learning algorithm returns a hypothesis that is consistent with this many examples, then with probability at least  $1 - \delta$ , it has error at most  $\epsilon$ . In other words, it is probably approximately correct. The number of required examples, as a function of  $\epsilon$  and  $\delta$ , is called the **sample complexity** of the hypothesis space.

SAMPLE  
COMPLEXITY

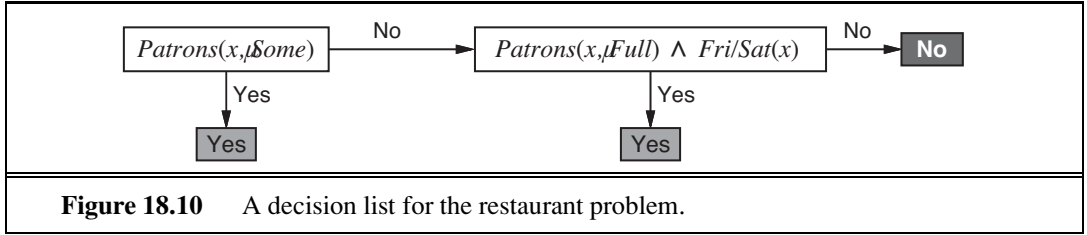
As we saw earlier, if  $\mathcal{H}$  is the set of all Boolean functions on  $n$  attributes, then  $|\mathcal{H}| = 2^{2^n}$ . Thus, the sample complexity of the space grows as  $2^n$ . Because the number of possible examples is also  $2^n$ , this suggests that PAC-learning in the class of all Boolean functions requires seeing all, or nearly all, of the possible examples. A moment's thought reveals the reason for this:  $\mathcal{H}$  contains enough hypotheses to classify any given set of examples in all possible ways. In particular, for any set of  $N$  examples, the set of hypotheses consistent with those examples contains equal numbers of hypotheses that predict  $x_{N+1}$  to be positive and hypotheses that predict  $x_{N+1}$  to be negative.

To obtain real generalization to unseen examples, then, it seems we need to restrict the hypothesis space  $\mathcal{H}$  in some way; but of course, if we do restrict the space, we might eliminate the true function altogether. There are three ways to escape this dilemma. The first, which we will cover in Chapter 19, is to bring prior knowledge to bear on the problem. The second, which we introduced in Section 18.4.3, is to insist that the algorithm return not just any consistent hypothesis, but preferably a simple one (as is done in decision tree learning). In cases where finding simple consistent hypotheses is tractable, the sample complexity results are generally better than for analyses based only on consistency. The third escape, which we pursue next, is to focus on learnable subsets of the entire hypothesis space of Boolean functions. This approach relies on the assumption that the restricted language contains a hypothesis  $h$  that is close enough to the true function  $f$ ; the benefits are that the restricted hypothesis space allows for effective generalization and is typically easier to search. We now examine one such restricted language in more detail.

### 18.5.1 PAC learning example: Learning decision lists

DECISION LISTS

We now show how to apply PAC learning to a new hypothesis space: **decision lists**. A decision list consists of a series of tests, each of which is a conjunction of literals. If a test succeeds when applied to an example description, the decision list specifies the value to be returned. If the test fails, processing continues with the next test in the list. Decision lists resemble decision trees, but their overall structure is simpler: they branch only in one



direction. In contrast, the individual tests are more complex. Figure 18.10 shows a decision list that represents the following hypothesis:

$$WillWait \Leftrightarrow (Patrons = Some) \vee (Patrons = Full \wedge Fri/Sat).$$

If we allow tests of arbitrary size, then decision lists can represent any Boolean function (Exercise 18.14). On the other hand, if we restrict the size of each test to at most  $k$  literals, then it is possible for the learning algorithm to generalize successfully from a small number of examples. We call this language  $k$ -DL. The example in Figure 18.10 is in 2-DL. It is easy to show (Exercise 18.14) that  $k$ -DL includes as a subset the language  $k$ -DT, the set of all decision trees of depth at most  $k$ . It is important to remember that the particular language referred to by  $k$ -DL depends on the attributes used to describe the examples. We will use the notation  $k$ -DL( $n$ ) to denote a  $k$ -DL language using  $n$  Boolean attributes.

The first task is to show that  $k$ -DL is learnable—that is, that any function in  $k$ -DL can be approximated accurately after training on a reasonable number of examples. To do this, we need to calculate the number of hypotheses in the language. Let the language of tests—conjunctions of at most  $k$  literals using  $n$  attributes—be  $Conj(n, k)$ . Because a decision list is constructed of tests, and because each test can be attached to either a *Yes* or a *No* outcome or can be absent from the decision list, there are at most  $3^{|Conj(n, k)|}$  distinct sets of component tests. Each of these sets of tests can be in any order, so

$$|k\text{-DL}(n)| \leq 3^{|Conj(n, k)|} |Conj(n, k)|!.$$

The number of conjunctions of  $k$  literals from  $n$  attributes is given by

$$|Conj(n, k)| = \sum_{i=0}^k \binom{2n}{i} = O(n^k).$$

Hence, after some work, we obtain

$$|k\text{-DL}(n)| = 2^{O(n^k \log_2(n^k))}.$$

We can plug this into Equation (18.1) to show that the number of examples needed for PAC-learning a  $k$ -DL function is polynomial in  $n$ :

$$N \geq \frac{1}{\epsilon} \left( \ln \frac{1}{\delta} + O(n^k \log_2(n^k)) \right).$$

Therefore, any algorithm that returns a consistent decision list will PAC-learn a  $k$ -DL function in a reasonable number of examples, for small  $k$ .

The next task is to find an efficient algorithm that returns a consistent decision list. We will use a greedy algorithm called DECISION-LIST-LEARNING that repeatedly finds a



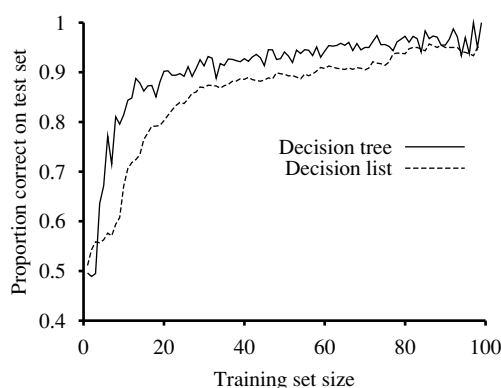
```

function DECISION-LIST-LEARNING(examples) returns a decision list, or failure

  if examples is empty then return the trivial decision list No
  t  $\leftarrow$  a test that matches a nonempty subset examplest of examples
    such that the members of examplest are all positive or all negative
  if there is no such t then return failure
  if the examples in examplest are positive then o  $\leftarrow$  Yes else o  $\leftarrow$  No
  return a decision list with initial test t and outcome o and remaining tests given by
    DECISION-LIST-LEARNING(examples − examplest)

```

**Figure 18.11** An algorithm for learning decision lists.



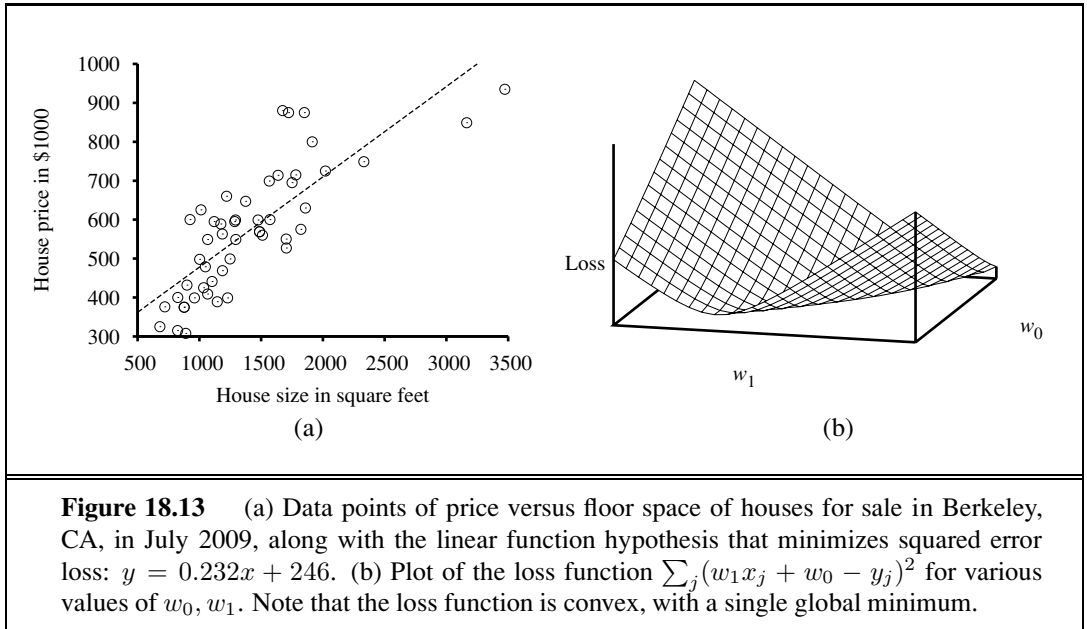
**Figure 18.12** Learning curve for DECISION-LIST-LEARNING algorithm on the restaurant data. The curve for DECISION-TREE-LEARNING is shown for comparison.

test that agrees exactly with some subset of the training set. Once it finds such a test, it adds it to the decision list under construction and removes the corresponding examples. It then constructs the remainder of the decision list, using just the remaining examples. This is repeated until there are no examples left. The algorithm is shown in Figure 18.11.

This algorithm does not specify the method for selecting the next test to add to the decision list. Although the formal results given earlier do not depend on the selection method, it would seem reasonable to prefer small tests that match large sets of uniformly classified examples, so that the overall decision list will be as compact as possible. The simplest strategy is to find the smallest test *t* that matches any uniformly classified subset, regardless of the size of the subset. Even this approach works quite well, as Figure 18.12 suggests.

## 18.6 REGRESSION AND CLASSIFICATION WITH LINEAR MODELS

Now it is time to move on from decision trees and lists to a different hypothesis space, one that has been used for hundred of years: the class of **linear functions** of continuous-valued



**Figure 18.13** (a) Data points of price versus floor space of houses for sale in Berkeley, CA, in July 2009, along with the linear function hypothesis that minimizes squared error loss:  $y = 0.232x + 246$ . (b) Plot of the loss function  $\sum_j (w_1 x_j + w_0 - y_j)^2$  for various values of  $w_0, w_1$ . Note that the loss function is convex, with a single global minimum.

inputs. We'll start with the simplest case: regression with a univariate linear function, otherwise known as “fitting a straight line.” Section 18.6.2 covers the multivariate case. Sections 18.6.3 and 18.6.4 show how to turn linear functions into classifiers by applying hard and soft thresholds.

### 18.6.1 Univariate linear regression

A univariate linear function (a straight line) with input  $x$  and output  $y$  has the form  $y = w_1 x + w_0$ , where  $w_0$  and  $w_1$  are real-valued coefficients to be learned. We use the letter  $w$  because we think of the coefficients as **weights**; the value of  $y$  is changed by changing the relative weight of one term or another. We'll define  $\mathbf{w}$  to be the vector  $[w_0, w_1]$ , and define

$$h_{\mathbf{w}}(x) = w_1 x + w_0.$$

Figure 18.13(a) shows an example of a training set of  $n$  points in the  $x, y$  plane, each point representing the size in square feet and the price of a house offered for sale. The task of finding the  $h_{\mathbf{w}}$  that best fits these data is called **linear regression**. To fit a line to the data, all we have to do is find the values of the weights  $[w_0, w_1]$  that minimize the empirical loss. It is traditional (going back to Gauss<sup>3</sup>) to use the squared loss function,  $L_2$ , summed over all the training examples:

$$\text{Loss}(h_{\mathbf{w}}) = \sum_{j=1}^N L_2(y_j, h_{\mathbf{w}}(x_j)) = \sum_{j=1}^N (y_j - h_{\mathbf{w}}(x_j))^2 = \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2.$$

<sup>3</sup> Gauss showed that if the  $y_j$  values have normally distributed noise, then the most likely values of  $w_1$  and  $w_0$  are obtained by minimizing the sum of the squares of the errors.

We would like to find  $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \operatorname{Loss}(h_{\mathbf{w}})$ . The sum  $\sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2$  is minimized when its partial derivatives with respect to  $w_0$  and  $w_1$  are zero:

$$\frac{\partial}{\partial w_0} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0 \text{ and } \frac{\partial}{\partial w_1} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0. \quad (18.2)$$

These equations have a unique solution:

$$w_1 = \frac{N(\sum x_j y_j) - (\sum x_j)(\sum y_j)}{N(\sum x_j^2) - (\sum x_j)^2}; \quad w_0 = (\sum y_j - w_1(\sum x_j))/N. \quad (18.3)$$

For the example in Figure 18.13(a), the solution is  $w_1 = 0.232$ ,  $w_0 = 246$ , and the line with those weights is shown as a dashed line in the figure.

WEIGHT SPACE

Many forms of learning involve adjusting weights to minimize a loss, so it helps to have a mental picture of what's going on in **weight space**—the space defined by all possible settings of the weights. For univariate linear regression, the weight space defined by  $w_0$  and  $w_1$  is two-dimensional, so we can graph the loss as a function of  $w_0$  and  $w_1$  in a 3D plot (see Figure 18.13(b)). We see that the loss function is **convex**, as defined on page 133; this is true for *every* linear regression problem with an  $L_2$  loss function, and implies that there are no local minima. In some sense that's the end of the story for linear models; if we need to fit lines to data, we apply Equation (18.3).<sup>4</sup>

GRADIENT DESCENT

To go beyond linear models, we will need to face the fact that the equations defining minimum loss (as in Equation (18.2)) will often have no closed-form solution. Instead, we will face a general optimization search problem in a continuous weight space. As indicated in Section 4.2 (page 129), such problems can be addressed by a hill-climbing algorithm that follows the **gradient** of the function to be optimized. In this case, because we are trying to minimize the loss, we will use **gradient descent**. We choose any starting point in weight space—here, a point in the  $(w_0, w_1)$  plane—and then move to a neighboring point that is downhill, repeating until we converge on the minimum possible loss:

$\mathbf{w} \leftarrow$  any point in the parameter space

**loop** until convergence **do**

**for each**  $w_i$  **in**  $\mathbf{w}$  **do**

$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \operatorname{Loss}(\mathbf{w}) \quad (18.4)$$

LEARNING RATE

The parameter  $\alpha$ , which we called the **step size** in Section 4.2, is usually called the **learning rate** when we are trying to minimize loss in a learning problem. It can be a fixed constant, or it can decay over time as the learning process proceeds.

For univariate regression, the loss function is a quadratic function, so the partial derivative will be a linear function. (The only calculus you need to know is that  $\frac{\partial}{\partial x} x^2 = 2x$  and  $\frac{\partial}{\partial x} x = 1$ .) Let's first work out the partial derivatives—the slopes—in the simplified case of

<sup>4</sup> With some caveats: the  $L_2$  loss function is appropriate when there is normally-distributed noise that is independent of  $x$ ; all results rely on the stationarity assumption; etc.

only one training example,  $(x, y)$ :

$$\begin{aligned}\frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w}) &= \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(x))^2 \\ &= 2(y - h_{\mathbf{w}}(x)) \times \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(x)) \\ &= 2(y - h_{\mathbf{w}}(x)) \times \frac{\partial}{\partial w_i} (y - (w_1 x + w_0)) ,\end{aligned}\tag{18.5}$$

applying this to both  $w_0$  and  $w_1$  we get:

$$\frac{\partial}{\partial w_0} \text{Loss}(\mathbf{w}) = -2(y - h_{\mathbf{w}}(x)); \quad \frac{\partial}{\partial w_1} \text{Loss}(\mathbf{w}) = -2(y - h_{\mathbf{w}}(x)) \times x$$

Then, plugging this back into Equation (18.4), and folding the 2 into the unspecified learning rate  $\alpha$ , we get the following learning rule for the weights:

$$w_0 \leftarrow w_0 + \alpha (y - h_{\mathbf{w}}(x)); \quad w_1 \leftarrow w_1 + \alpha (y - h_{\mathbf{w}}(x)) \times x$$

These updates make intuitive sense: if  $h_{\mathbf{w}}(x) > y$ , i.e., the output of the hypothesis is too large, reduce  $w_0$  a bit, and reduce  $w_1$  if  $x$  was a positive input but increase  $w_1$  if  $x$  was a negative input.

The preceding equations cover one training example. For  $N$  training examples, we want to minimize the sum of the individual losses for each example. The derivative of a sum is the sum of the derivatives, so we have:

$$w_0 \leftarrow w_0 + \alpha \sum_j (y_j - h_{\mathbf{w}}(x_j)); \quad w_1 \leftarrow w_1 + \alpha \sum_j (y_j - h_{\mathbf{w}}(x_j)) \times x_j .$$

BATCH GRADIENT  
DESCENT

These updates constitute the **batch gradient descent** learning rule for univariate linear regression. Convergence to the unique global minimum is guaranteed (as long as we pick  $\alpha$  small enough) but may be very slow: we have to cycle through all the training data for every step, and there may be many steps.

STOCHASTIC  
GRADIENT DESCENT

There is another possibility, called **stochastic gradient descent**, where we consider only a single training point at a time, taking a step after each one using Equation (18.5). Stochastic gradient descent can be used in an online setting, where new data are coming in one at a time, or offline, where we cycle through the same data as many times as is necessary, taking a step after considering each single example. It is often faster than batch gradient descent. With a fixed learning rate  $\alpha$ , however, it does not guarantee convergence; it can oscillate around the minimum without settling down. In some cases, as we see later, a schedule of decreasing learning rates (as in simulated annealing) does guarantee convergence.

## 18.6.2 Multivariate linear regression

MULTIVARIATE  
LINEAR REGRESSION

We can easily extend to **multivariate linear regression** problems, in which each example  $\mathbf{x}_j$  is an  $n$ -element vector.<sup>5</sup> Our hypothesis space is the set of functions of the form

$$h_{sw}(\mathbf{x}_j) = w_0 + w_1 x_{j,1} + \cdots + w_n x_{j,n} = w_0 + \sum_i w_i x_{j,i} .$$

<sup>5</sup> The reader may wish to consult Appendix A for a brief summary of linear algebra.

The  $w_0$  term, the intercept, stands out as different from the others. We can fix that by inventing a dummy input attribute,  $x_{j,0}$ , which is defined as always equal to 1. Then  $h$  is simply the dot product of the weights and the input vector (or equivalently, the matrix product of the transpose of the weights and the input vector):

$$h_{sw}(\mathbf{x}_j) = \mathbf{w} \cdot \mathbf{x}_j = \mathbf{w}^\top \mathbf{x}_j = \sum_i w_i x_{j,i}.$$

The best vector of weights,  $\mathbf{w}^*$ , minimizes squared-error loss over the examples:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_j L_2(y_j, \mathbf{w} \cdot \mathbf{x}_j).$$

Multivariate linear regression is actually not much more complicated than the univariate case we just covered. Gradient descent will reach the (unique) minimum of the loss function; the update equation for each weight  $w_i$  is

$$w_i \leftarrow w_i + \alpha \sum_j x_{j,i}(y_j - h_{\mathbf{w}}(\mathbf{x}_j)). \quad (18.6)$$

It is also possible to solve analytically for the  $\mathbf{w}$  that minimizes loss. Let  $\mathbf{y}$  be the vector of outputs for the training examples, and  $\mathbf{X}$  be the **data matrix**, i.e., the matrix of inputs with one  $n$ -dimensional example per row. Then the solution

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

minimizes the squared error.

With univariate linear regression we didn't have to worry about overfitting. But with multivariate linear regression in high-dimensional spaces it is possible that some dimension that is actually irrelevant appears by chance to be useful, resulting in **overfitting**.

Thus, it is common to use **regularization** on multivariate linear functions to avoid overfitting. Recall that with regularization we minimize the total cost of a hypothesis, counting both the empirical loss and the complexity of the hypothesis:

$$\operatorname{Cost}(h) = \operatorname{EmpLoss}(h) + \lambda \operatorname{Complexity}(h).$$

For linear functions the complexity can be specified as a function of the weights. We can consider a family of regularization functions:

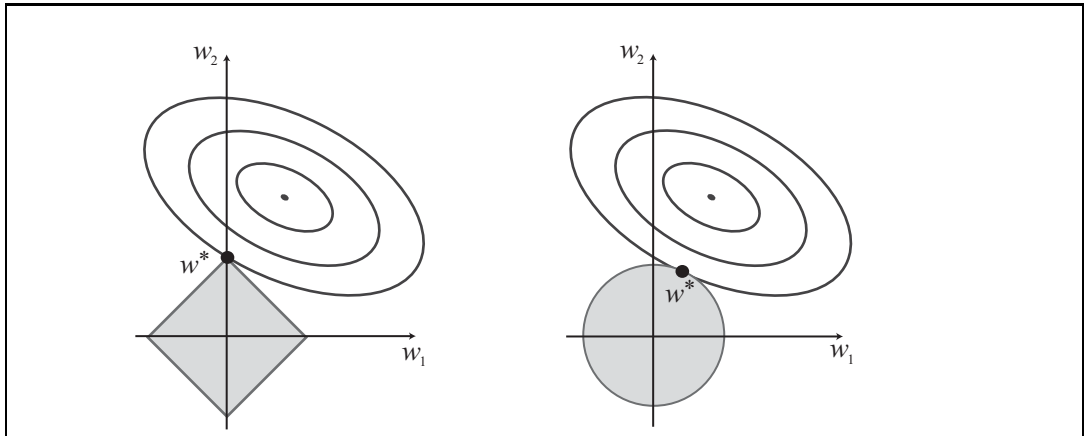
$$\operatorname{Complexity}(h_{\mathbf{w}}) = L_q(\mathbf{w}) = \sum_i |w_i|^q.$$

As with loss functions,<sup>6</sup> with  $q=1$  we have  $L_1$  regularization, which minimizes the sum of the absolute values; with  $q=2$ ,  $L_2$  regularization minimizes the sum of squares. Which regularization function should you pick? That depends on the specific problem, but  $L_1$  regularization has an important advantage: it tends to produce a **sparse model**. That is, it often sets many weights to zero, effectively declaring the corresponding attributes to be irrelevant—just as DECISION-TREE-LEARNING does (although by a different mechanism). Hypotheses that discard attributes can be easier for a human to understand, and may be less likely to overfit.

<sup>6</sup> It is perhaps confusing that  $L_1$  and  $L_2$  are used for both loss functions and regularization functions. They need not be used in pairs: you could use  $L_2$  loss with  $L_1$  regularization, or vice versa.

DATA MATRIX

SPARSE MODEL



**Figure 18.14** Why  $L_1$  regularization tends to produce a sparse model. (a) With  $L_1$  regularization (box), the minimal achievable loss (concentric contours) often occurs on an axis, meaning a weight of zero. (b) With  $L_2$  regularization (circle), the minimal loss is likely to occur anywhere on the circle, giving no preference to zero weights.

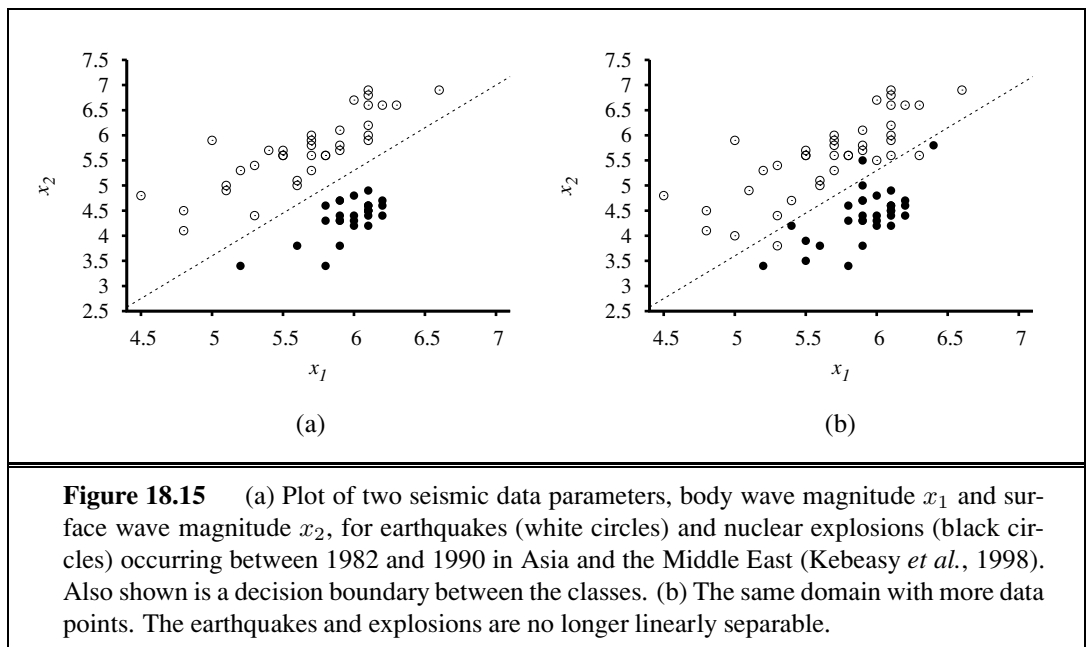
Figure 18.14 gives an intuitive explanation of why  $L_1$  regularization leads to weights of zero, while  $L_2$  regularization does not. Note that minimizing  $Loss(\mathbf{w}) + \lambda Complexity(\mathbf{w})$  is equivalent to minimizing  $Loss(\mathbf{w})$  subject to the constraint that  $Complexity(\mathbf{w}) \leq c$ , for some constant  $c$  that is related to  $\lambda$ . Now, in Figure 18.14(a) the diamond-shaped box represents the set of points  $\mathbf{w}$  in two-dimensional weight space that have  $L_1$  complexity less than  $c$ ; our solution will have to be somewhere inside this box. The concentric ovals represent contours of the loss function, with the minimum loss at the center. We want to find the point in the box that is closest to the minimum; you can see from the diagram that, for an arbitrary position of the minimum and its contours, it will be common for the corner of the box to find its way closest to the minimum, just because the corners are pointy. And of course the corners are the points that have a value of zero in some dimension. In Figure 18.14(b), we've done the same for the  $L_2$  complexity measure, which represents a circle rather than a diamond. Here you can see that, in general, there is no reason for the intersection to appear on one of the axes; thus  $L_2$  regularization does not tend to produce zero weights. The result is that the number of examples required to find a good  $h$  is linear in the number of irrelevant features for  $L_2$  regularization, but only logarithmic with  $L_1$  regularization. Empirical evidence on many problems supports this analysis.

Another way to look at it is that  $L_1$  regularization takes the dimensional axes seriously, while  $L_2$  treats them as arbitrary. The  $L_2$  function is spherical, which makes it rotationally invariant: Imagine a set of points in a plane, measured by their  $x$  and  $y$  coordinates. Now imagine rotating the axes by  $45^\circ$ . You'd get a different set of  $(x', y')$  values representing the same points. If you apply  $L_2$  regularization before and after rotating, you get exactly the same point as the answer (although the point would be described with the new  $(x', y')$  coordinates). That is appropriate when the choice of axes really is arbitrary—when it doesn't matter whether your two dimensions are distances north and east; or distances north-east and

south-east. With  $L_1$  regularization you'd get a different answer, because the  $L_1$  function is not rotationally invariant. That is appropriate when the axes are not interchangeable; it doesn't make sense to rotate "number of bathrooms"  $45^\circ$  towards "lot size."

### 18.6.3 Linear classifiers with a hard threshold

Linear functions can be used to do classification as well as regression. For example, Figure 18.15(a) shows data points of two classes: earthquakes (which are of interest to seismologists) and underground explosions (which are of interest to arms control experts). Each point is defined by two input values,  $x_1$  and  $x_2$ , that refer to body and surface wave magnitudes computed from the seismic signal. Given these training data, the task of classification is to learn a hypothesis  $h$  that will take new  $(x_1, x_2)$  points and return either 0 for earthquakes or 1 for explosions.



**Figure 18.15** (a) Plot of two seismic data parameters, body wave magnitude  $x_1$  and surface wave magnitude  $x_2$ , for earthquakes (white circles) and nuclear explosions (black circles) occurring between 1982 and 1990 in Asia and the Middle East (Kebeasy *et al.*, 1998). Also shown is a decision boundary between the classes. (b) The same domain with more data points. The earthquakes and explosions are no longer linearly separable.

DECISION  
BOUNDARY

LINEAR SEPARATOR  
LINEAR  
SEPARABILITY

A **decision boundary** is a line (or a surface, in higher dimensions) that separates the two classes. In Figure 18.15(a), the decision boundary is a straight line. A linear decision boundary is called a **linear separator** and data that admit such a separator are called **linearly separable**. The linear separator in this case is defined by

$$x_2 = 1.7x_1 - 4.9 \quad \text{or} \quad -4.9 + 1.7x_1 - x_2 = 0.$$

The explosions, which we want to classify with value 1, are to the right of this line with higher values of  $x_1$  and lower values of  $x_2$ , so they are points for which  $-4.9 + 1.7x_1 - x_2 > 0$ , while earthquakes have  $-4.9 + 1.7x_1 - x_2 < 0$ . Using the convention of a dummy input  $x_0 = 1$ , we can write the classification hypothesis as

$$h_{\mathbf{w}}(\mathbf{x}) = 1 \text{ if } \mathbf{w} \cdot \mathbf{x} \geq 0 \text{ and } 0 \text{ otherwise.}$$

THRESHOLD  
FUNCTION

Alternatively, we can think of  $h$  as the result of passing the linear function  $\mathbf{w} \cdot \mathbf{x}$  through a **threshold function**:

$$h_{\mathbf{w}}(\mathbf{x}) = \text{Threshold}(\mathbf{w} \cdot \mathbf{x}) \text{ where } \text{Threshold}(z) = 1 \text{ if } z \geq 0 \text{ and } 0 \text{ otherwise.}$$

The threshold function is shown in Figure 18.17(a).

Now that the hypothesis  $h_{\mathbf{w}}(\mathbf{x})$  has a well-defined mathematical form, we can think about choosing the weights  $\mathbf{w}$  to minimize the loss. In Sections 18.6.1 and 18.6.2, we did this both in closed form (by setting the gradient to zero and solving for the weights) and by gradient descent in weight space. Here, we cannot do either of those things because the gradient is zero almost everywhere in weight space except at those points where  $\mathbf{w} \cdot \mathbf{x} = 0$ , and at those points the gradient is undefined.

There is, however, a simple weight update rule that converges to a solution—that is, a linear separator that classifies the data perfectly—provided the data are linearly separable. For a single example  $(\mathbf{x}, y)$ , we have

$$w_i \leftarrow w_i + \alpha (y - h_{\mathbf{w}}(\mathbf{x})) \times x_i \quad (18.7)$$

PERCEPTRON  
LEARNING RULE

which is essentially identical to the Equation (18.6), the update rule for linear regression! This rule is called the **perceptron learning rule**, for reasons that will become clear in Section 18.7. Because we are considering a 0/1 classification problem, however, the behavior is somewhat different. Both the true value  $y$  and the hypothesis output  $h_{\mathbf{w}}(\mathbf{x})$  are either 0 or 1, so there are three possibilities:

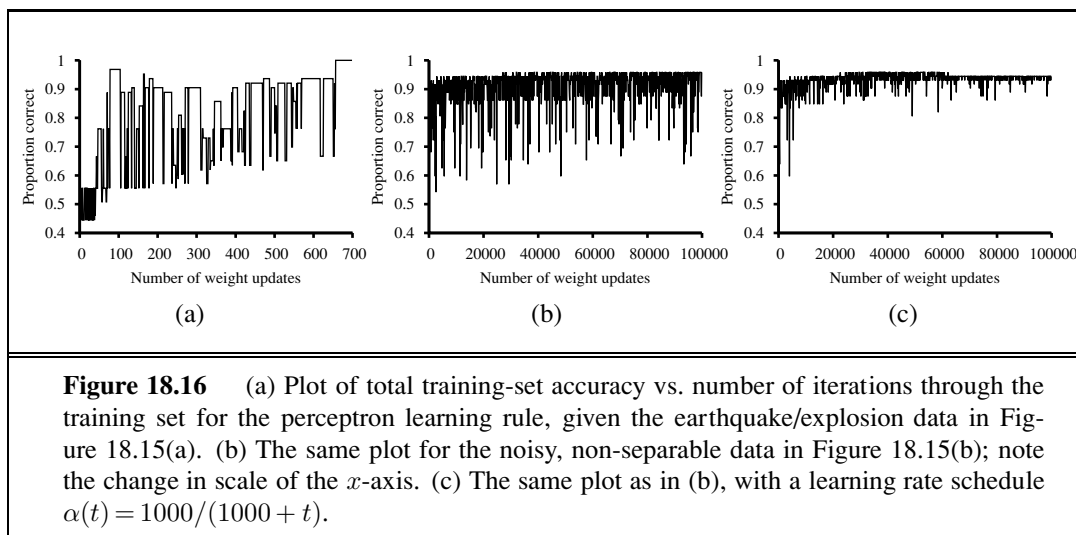
- If the output is correct, i.e.,  $y = h_{\mathbf{w}}(\mathbf{x})$ , then the weights are not changed.
- If  $y$  is 1 but  $h_{\mathbf{w}}(\mathbf{x})$  is 0, then  $w_i$  is *increased* when the corresponding input  $x_i$  is positive and *decreased* when  $x_i$  is negative. This makes sense, because we want to make  $\mathbf{w} \cdot \mathbf{x}$  bigger so that  $h_{\mathbf{w}}(\mathbf{x})$  outputs a 1.
- If  $y$  is 0 but  $h_{\mathbf{w}}(\mathbf{x})$  is 1, then  $w_i$  is *decreased* when the corresponding input  $x_i$  is positive and *increased* when  $x_i$  is negative. This makes sense, because we want to make  $\mathbf{w} \cdot \mathbf{x}$  smaller so that  $h_{\mathbf{w}}(\mathbf{x})$  outputs a 0.

## TRAINING CURVE

Typically the learning rule is applied one example at a time, choosing examples at random (as in stochastic gradient descent). Figure 18.16(a) shows a **training curve** for this learning rule applied to the earthquake/explosion data shown in Figure 18.15(a). A training curve measures the classifier performance on a fixed training set as the learning process proceeds on that same training set. The curve shows the update rule converging to a zero-error linear separator. The “convergence” process isn’t exactly pretty, but it always works. This particular run takes 657 steps to converge, for a data set with 63 examples, so each example is presented roughly 10 times on average. Typically, the variation across runs is very large.

We have said that the perceptron learning rule converges to a perfect linear separator when the data points are linearly separable, but what if they are not? This situation is all too common in the real world. For example, Figure 18.15(b) adds back in the data points left out by Kebeasy *et al.* (1998) when they plotted the data shown in Figure 18.15(a). In Figure 18.16(b), we show the perceptron learning rule failing to converge even after 10,000 steps: even though it hits the minimum-error solution (three errors) many times, the algorithm keeps changing the weights. In general, the perceptron rule may not converge to a





stable solution for fixed learning rate  $\alpha$ , but if  $\alpha$  decays as  $O(1/t)$  where  $t$  is the iteration number, then the rule can be shown to converge to a minimum-error solution when examples are presented in a random sequence.<sup>7</sup> It can also be shown that finding the minimum-error solution is NP-hard, so one expects that many presentations of the examples will be required for convergence to be achieved. Figure 18.16(b) shows the training process with a learning rate schedule  $\alpha(t) = 1000/(1000 + t)$ : convergence is not perfect after 100,000 iterations, but it is much better than the fixed- $\alpha$  case.

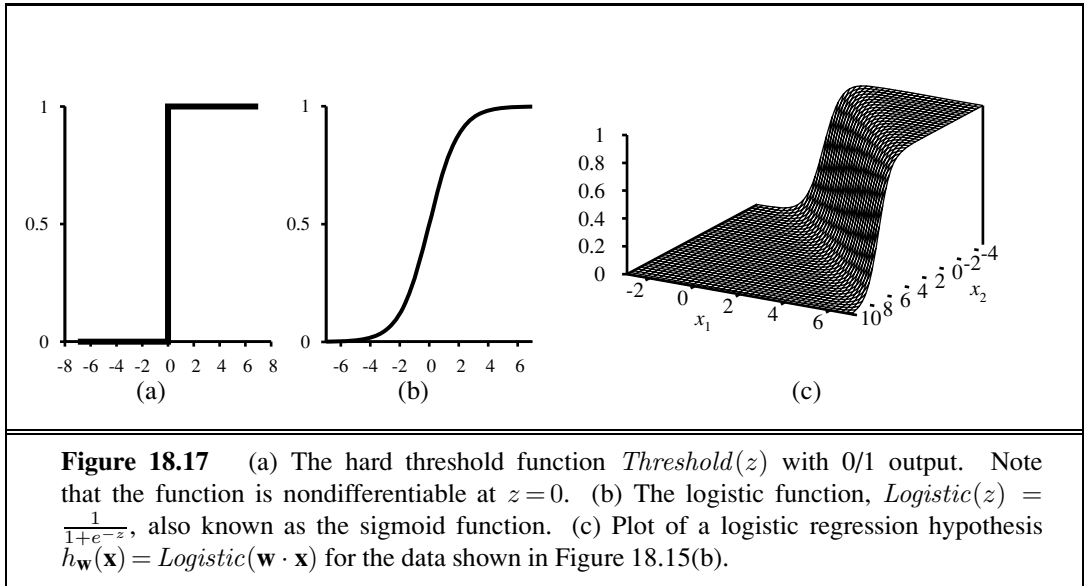
#### 18.6.4 Linear classification with logistic regression

We have seen that passing the output of a linear function through the threshold function creates a linear classifier; yet the hard nature of the threshold causes some problems: the hypothesis  $h_{\mathbf{w}}(\mathbf{x})$  is not differentiable and is in fact a discontinuous function of its inputs and its weights; this makes learning with the perceptron rule a very unpredictable adventure. Furthermore, the linear classifier always announces a completely confident prediction of 1 or 0, even for examples that are very close to the boundary; in many situations, we really need more gradated predictions.

All of these issues can be resolved to a large extent by softening the threshold function—approximating the hard threshold with a continuous, differentiable function. In Chapter 14 (page 522), we saw two functions that look like soft thresholds: the integral of the standard normal distribution (used for the probit model) and the logistic function (used for the logit model). Although the two functions are very similar in shape, the logistic function

$$\text{Logistic}(z) = \frac{1}{1 + e^{-z}}$$

<sup>7</sup> Technically, we require that  $\sum_{t=1}^{\infty} \alpha(t) = \infty$  and  $\sum_{t=1}^{\infty} \alpha^2(t) < \infty$ . The decay  $\alpha(t) = O(1/t)$  satisfies these conditions.



has more convenient mathematical properties. The function is shown in Figure 18.17(b). With the logistic function replacing the threshold function, we now have

$$h_{\mathbf{w}}(\mathbf{x}) = Logistic(\mathbf{w} \cdot \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}.$$

An example of such a hypothesis for the two-input earthquake/explosion problem is shown in Figure 18.17(c). Notice that the output, being a number between 0 and 1, can be interpreted as a *probability* of belonging to the class labeled 1. The hypothesis forms a soft boundary in the input space and gives a probability of 0.5 for any input at the center of the boundary region, and approaches 0 or 1 as we move away from the boundary.

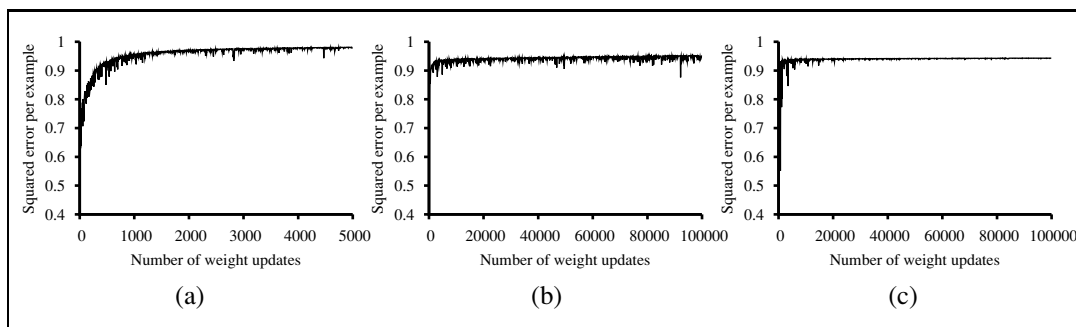
The process of fitting the weights of this model to minimize loss on a data set is called **logistic regression**. There is no easy closed-form solution to find the optimal value of  $\mathbf{w}$  with this model, but the gradient descent computation is straightforward. Because our hypotheses no longer output just 0 or 1, we will use the  $L_2$  loss function; also, to keep the formulas readable, we'll use  $g$  to stand for the logistic function, with  $g'$  its derivative.

For a single example  $(\mathbf{x}, y)$ , the derivation of the gradient is the same as for linear regression (Equation (18.5)) up to the point where the actual form of  $h$  is inserted. (For this derivation, we will need the **chain rule**:  $\partial g(f(x))/\partial x = g'(f(x)) \partial f(x)/\partial x$ .) We have

$$\begin{aligned} \frac{\partial}{\partial w_i} Loss(\mathbf{w}) &= \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(\mathbf{x}))^2 \\ &= 2(y - h_{\mathbf{w}}(\mathbf{x})) \times \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(\mathbf{x})) \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \times g'(\mathbf{w} \cdot \mathbf{x}) \times \frac{\partial}{\partial w_i} \mathbf{w} \cdot \mathbf{x} \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \times g'(\mathbf{w} \cdot \mathbf{x}) \times x_i. \end{aligned}$$

LOGISTIC  
REGRESSION

CHAIN RULE



**Figure 18.18** Repeat of the experiments in Figure 18.16 using logistic regression and squared error. The plot in (a) covers 5000 iterations rather than 1000, while (b) and (c) use the same scale.

The derivative  $g'$  of the logistic function satisfies  $g'(z) = g(z)(1 - g(z))$ , so we have

$$g'(\mathbf{w} \cdot \mathbf{x}) = g(\mathbf{w} \cdot \mathbf{x})(1 - g(\mathbf{w} \cdot \mathbf{x})) = h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x}))$$

so the weight update for minimizing the loss is

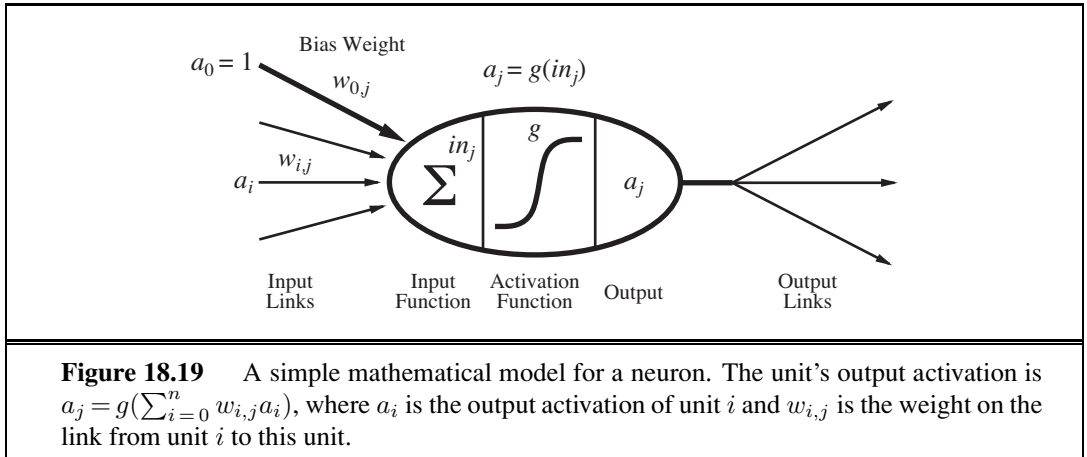
$$w_i \leftarrow w_i + \alpha (y - h_{\mathbf{w}}(\mathbf{x})) \times h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x})) \times x_i. \quad (18.8)$$

Repeating the experiments of Figure 18.16 with logistic regression instead of the linear threshold classifier, we obtain the results shown in Figure 18.18. In (a), the linearly separable case, logistic regression is somewhat slower to converge, but behaves much more predictably. In (b) and (c), where the data are noisy and nonseparable, logistic regression converges far more quickly and reliably. These advantages tend to carry over into real-world applications and logistic regression has become one of the most popular classification techniques for problems in medicine, marketing and survey analysis, credit scoring, public health, and other applications.

## 18.7 ARTIFICIAL NEURAL NETWORKS

We turn now to what seems to be a somewhat unrelated topic: the brain. In fact, as we will see, the technical ideas we have discussed so far in this chapter turn out to be useful in building mathematical models of the brain's activity; conversely, thinking about the brain has helped in extending the scope of the technical ideas.

Chapter 1 touched briefly on the basic findings of neuroscience—in particular, the hypothesis that mental activity consists primarily of electrochemical activity in networks of brain cells called **neurons**. (Figure 1.2 on page 11 showed a schematic diagram of a typical neuron.) Inspired by this hypothesis, some of the earliest AI work aimed to create artificial **neural networks**. (Other names for the field include **connectionism**, **parallel distributed processing**, and **neural computation**.) Figure 18.19 shows a simple mathematical model of the neuron devised by McCulloch and Pitts (1943). Roughly speaking, it “fires” when a linear combination of its inputs exceeds some (hard or soft) threshold—that is, it implements



a linear classifier of the kind described in the preceding section. A neural network is just a collection of units connected together; the properties of the network are determined by its topology and the properties of the “neurons.”

Since 1943, much more detailed and realistic models have been developed, both for neurons and for larger systems in the brain, leading to the modern field of **computational neuroscience**. On the other hand, researchers in AI and statistics became interested in the more abstract properties of neural networks, such as their ability to perform distributed computation, to tolerate noisy inputs, and to learn. Although we understand now that other kinds of systems—including Bayesian networks—have these properties, neural networks remain one of the most popular and effective forms of learning system and are worthy of study in their own right.

### 18.7.1 Neural network structures

Neural networks are composed of nodes or **units** (see Figure 18.19) connected by directed **links**. A link from unit  $i$  to unit  $j$  serves to propagate the **activation**  $a_i$  from  $i$  to  $j$ .<sup>8</sup> Each link also has a numeric **weight**  $w_{i,j}$  associated with it, which determines the strength and sign of the connection. Just as in linear regression models, each unit has a dummy input  $a_0 = 1$  with an associated weight  $w_{0,j}$ . Each unit  $j$  first computes a weighted sum of its inputs:

$$in_j = \sum_{i=0}^n w_{i,j} a_i .$$

Then it applies an **activation function**  $g$  to this sum to derive the output:

$$a_j = g(in_j) = g \left( \sum_{i=0}^n w_{i,j} a_i \right) . \quad (18.9)$$

<sup>8</sup> A note on notation: for this section, we are forced to suspend our usual conventions. Input attributes are still indexed by  $i$ , so that an “external” activation  $a_i$  is given by input  $x_i$ ; but index  $j$  will refer to internal units rather than examples. Throughout this section, the mathematical derivations concern a single generic example  $\mathbf{x}$ , omitting the usual summations over examples to obtain results for the whole data set.

PERCEPTRON  
SIGMOID  
PERCEPTRON

The activation function  $g$  is typically either a hard threshold (Figure 18.17(a)), in which case the unit is called a **perceptron**, or a logistic function (Figure 18.17(b)), in which case the term **sigmoid perceptron** is sometimes used. Both of these nonlinear activation function ensure the important property that the entire network of units can represent a nonlinear function (see Exercise 18.22). As mentioned in the discussion of logistic regression (page 725), the logistic activation function has the added advantage of being differentiable.

FEED-FORWARD  
NETWORK

RECURRENT  
NETWORK

Having decided on the mathematical model for individual “neurons,” the next task is to connect them together to form a network. There are two fundamentally distinct ways to do this. A **feed-forward network** has connections only in one direction—that is, it forms a directed acyclic graph. Every node receives input from “upstream” nodes and delivers output to “downstream” nodes; there are no loops. A feed-forward network represents a function of its current input; thus, it has no internal state other than the weights themselves. A **recurrent network**, on the other hand, feeds its outputs back into its own inputs. This means that the activation levels of the network form a dynamical system that may reach a stable state or exhibit oscillations or even chaotic behavior. Moreover, the response of the network to a given input depends on its initial state, which may depend on previous inputs. Hence, recurrent networks (unlike feed-forward networks) can support short-term memory. This makes them more interesting as models of the brain, but also more difficult to understand. This section will concentrate on feed-forward networks; some pointers for further reading on recurrent networks are given at the end of the chapter.

LAYERS

HIDDEN UNIT

Feed-forward networks are usually arranged in **layers**, such that each unit receives input only from units in the immediately preceding layer. In the next two subsections, we will look at single-layer networks, in which every unit connects directly from the network’s inputs to its outputs, and multilayer networks, which have one or more layers of **hidden units** that are not connected to the outputs of the network. So far in this chapter, we have considered only learning problems with a single output variable  $y$ , but neural networks are often used in cases where multiple outputs are appropriate. For example, if we want to train a network to add two input bits, each a 0 or a 1, we will need one output for the sum bit and one for the carry bit. Also, when the learning problem involves classification into more than two classes—for example, when learning to categorize images of handwritten digits—it is common to use one output unit for each class.

18.7.2 Single-layer feed-forward neural networks (perceptrons)

PERCEPTRON  
NETWORK

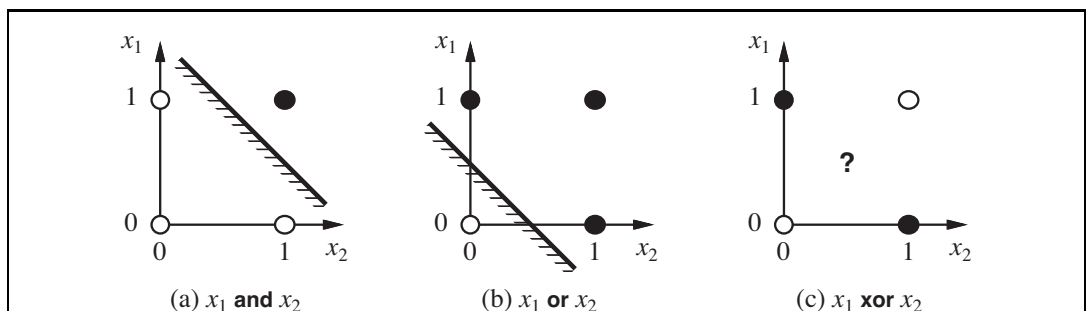
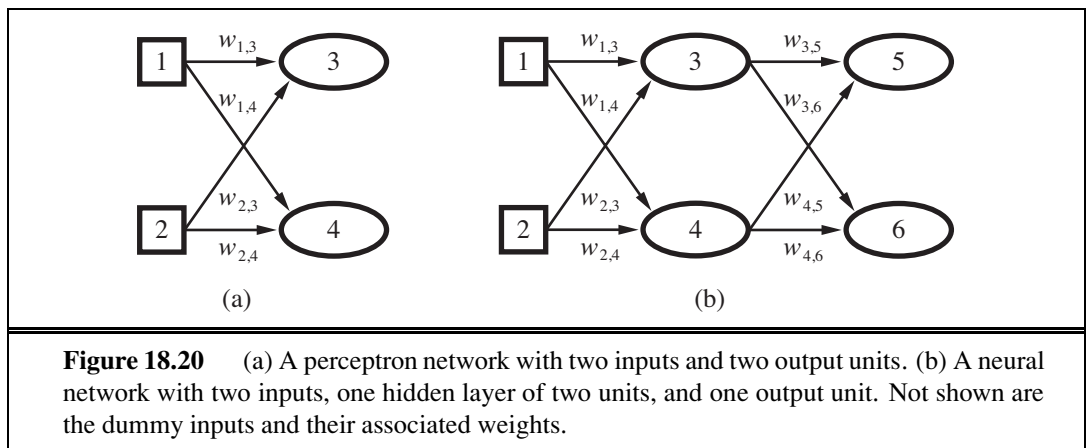
A network with all the inputs connected directly to the outputs is called a **single-layer neural network**, or a **perceptron network**. Figure 18.20 shows a simple two-input, two-output perceptron network. With such a network, we might hope to learn the two-bit adder function, for example. Here are all the training data we will need:

$x_1$	$x_2$	$y_3$ (carry)	$y_4$ (sum)
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0

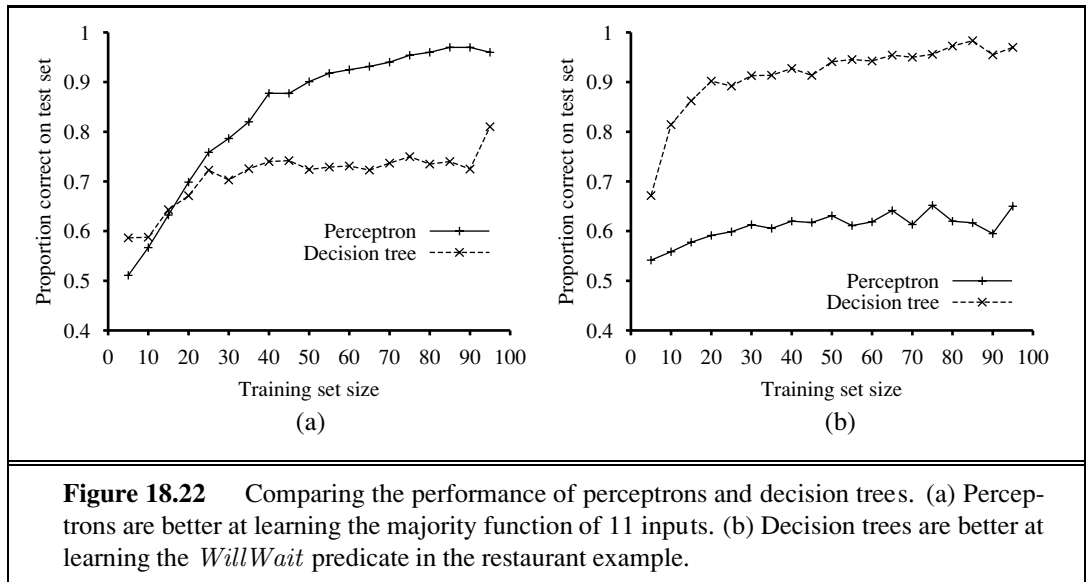
The first thing to notice is that a perceptron network with  $m$  outputs is really  $m$  separate networks, because each weight affects only one of the outputs. Thus, there will be  $m$  separate training processes. Furthermore, depending on the type of activation function used, the training processes will be either the **perceptron learning rule** (Equation (18.7) on page 724) or gradient descent rule for the **logistic regression** (Equation (18.8) on page 727).

If you try either method on the two-bit-adder data, something interesting happens. Unit 3 learns the carry function easily, but unit 4 completely fails to learn the sum function. No, unit 4 is not defective! The problem is with the sum function itself. We saw in Section 18.6 that linear classifiers (whether hard or soft) can represent linear decision boundaries in the input space. This works fine for the carry function, which is a logical AND (see Figure 18.21(a)). The sum function, however, is an XOR (exclusive OR) of the two inputs. As Figure 18.21(c) illustrates, this function is not linearly separable so the perceptron cannot learn it.

The linearly separable functions constitute just a small fraction of all Boolean functions; Exercise 18.20 asks you to quantify this fraction. The inability of perceptrons to learn even such simple functions as XOR was a significant setback to the nascent neural network



**Figure 18.21** Linear separability in threshold perceptrons. Black dots indicate a point in the input space where the value of the function is 1, and white dots indicate a point where the value is 0. The perceptron returns 1 on the region on the non-shaded side of the line. In (c), no such line exists that correctly classifies the inputs.



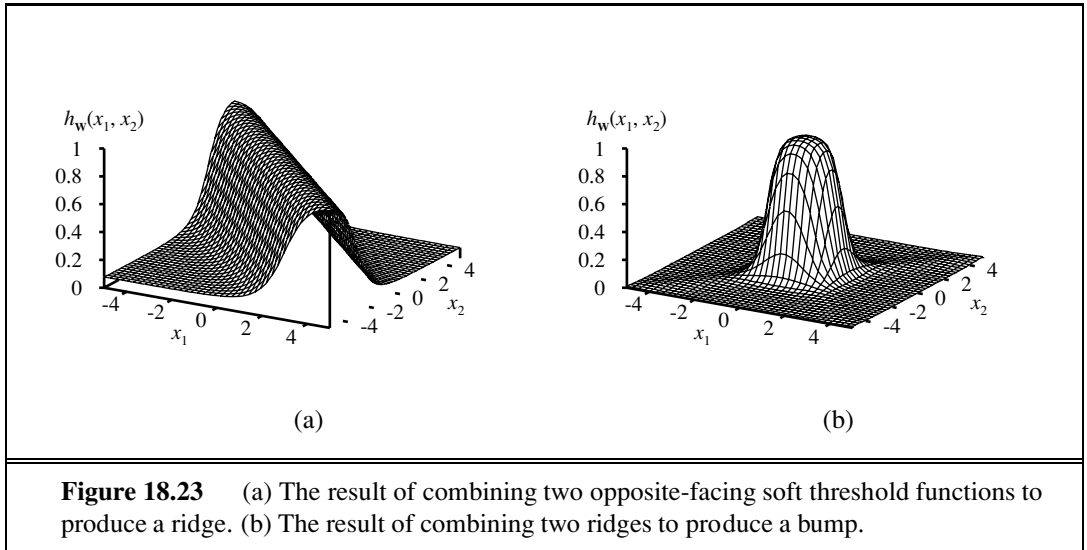
community in the 1960s. Perceptrons are far from useless, however. Section 18.6.4 noted that logistic regression (i.e., training a sigmoid perceptron) is even today a very popular and effective tool. Moreover, a perceptron can represent some quite “complex” Boolean functions very compactly. For example, the **majority function**, which outputs a 1 only if more than half of its  $n$  inputs are 1, can be represented by a perceptron with each  $w_i = 1$  and with  $w_0 = -n/2$ . A decision tree would need exponentially many nodes to represent this function.

Figure 18.22 shows the learning curve for a perceptron on two different problems. On the left, we show the curve for learning the majority function with 11 Boolean inputs (i.e., the function outputs a 1 if 6 or more inputs are 1). As we would expect, the perceptron learns the function quite quickly, because the majority function is linearly separable. On the other hand, the decision-tree learner makes no progress, because the majority function is very hard (although not impossible) to represent as a decision tree. On the right, we have the restaurant example. The solution problem is easily represented as a decision tree, but is not linearly separable. The best plane through the data correctly classifies only 65%.

### 18.7.3 Multilayer feed-forward neural networks

(McCulloch and Pitts, 1943) were well aware that a single threshold unit would not solve all their problems. In fact, their paper proves that such a unit can represent the basic Boolean functions AND, OR, and NOT and then goes on to argue that any desired functionality can be obtained by connecting large numbers of units into (possibly recurrent) networks of arbitrary depth. The problem was that nobody knew how to train such networks.

This turns out to be an easy problem if we think of a network the right way: as a function  $h_{\mathbf{w}}(\mathbf{x})$  parameterized by the weights  $\mathbf{w}$ . Consider the simple network shown in Figure 18.20(b), which has two input units, two hidden units, and two output unit. (In addition, each unit has a dummy input fixed at 1.) Given an input vector  $\mathbf{x} = (x_1, x_2)$ , the activations



of the input units are set to  $(a_1, a_2) = (x_1, x_2)$ . The output at unit 5 is given by

$$\begin{aligned}
 a_5 &= g(w_{0,5} + w_{3,5} a_3 + w_{4,5} a_4) \\
 &= g(w_{0,5} + w_{3,5} g(w_{0,3} + w_{1,3} a_1 + w_{2,3} a_2) + w_{4,5} g(w_{0,4} + w_{1,4} a_1 + w_{2,4} a_2)) \\
 &= g(w_{0,5} + w_{3,5} g(w_{0,3} + w_{1,3} x_1 + w_{2,3} x_2) + w_{4,5} g(w_{0,4} + w_{1,4} x_1 + w_{2,4} x_2)).
 \end{aligned}$$

Thus, we have the output expressed as a function of the inputs and the weights. A similar expression holds for unit 6. As long as we can calculate the derivatives of such expressions with respect to the weights, we can use the gradient-descent loss-minimization method to train the network. Section 18.7.4 shows exactly how to do this. And because the function represented by a network can be highly nonlinear—composed, as it is, of nested nonlinear soft threshold functions—we can see neural networks as a tool for doing **nonlinear regression**.

Before delving into learning rules, let us look at the ways in which networks generate complicated functions. First, remember that each unit in a sigmoid network represents a soft threshold in its input space, as shown in Figure 18.17(c) (page 726). With one hidden layer and one output layer, as in Figure 18.20(b), each output unit computes a soft-thresholded linear combination of several such functions. For example, by adding two opposite-facing soft threshold functions and thresholding the result, we can obtain a “ridge” function as shown in Figure 18.23(a). Combining two such ridges at right angles to each other (i.e., combining the outputs from four hidden units), we obtain a “bump” as shown in Figure 18.23(b).

With more hidden units, we can produce more bumps of different sizes in more places. In fact, with a single, sufficiently large hidden layer, it is possible to represent any continuous function of the inputs with arbitrary accuracy; with two layers, even discontinuous functions can be represented.<sup>9</sup> Unfortunately, for any *particular* network structure, it is harder to characterize exactly which functions can be represented and which ones cannot.

<sup>9</sup> The proof is complex, but the main point is that the required number of hidden units grows exponentially with the number of inputs. For example,  $2^n/n$  hidden units are needed to encode all Boolean functions of  $n$  inputs.



### 18.7.4 Learning in multilayer networks

First, let us dispense with one minor complication arising in multilayer networks: interactions among the learning problems when the network has multiple outputs. In such cases, we should think of the network as implementing a vector function  $\mathbf{h}_{\mathbf{w}}$  rather than a scalar function  $h_{\mathbf{w}}$ ; for example, the network in Figure 18.20(b) returns a vector  $[a_5, a_6]$ . Similarly, the target output will be a vector  $\mathbf{y}$ . Whereas a perceptron network decomposes into  $m$  separate learning problems for an  $m$ -output problem, this decomposition fails in a multilayer network. For example, both  $a_5$  and  $a_6$  in Figure 18.20(b) depend on all of the input-layer weights, so updates to those weights will depend on errors in both  $a_5$  and  $a_6$ . Fortunately, this dependency is very simple in the case of any loss function that is *additive* across the components of the error vector  $\mathbf{y} - \mathbf{h}_{\mathbf{w}}(\mathbf{x})$ . For the  $L_2$  loss, we have, for any weight  $w$ ,

$$\frac{\partial}{\partial w} \text{Loss}(\mathbf{w}) = \frac{\partial}{\partial w} |\mathbf{y} - \mathbf{h}_{\mathbf{w}}(\mathbf{x})|^2 = \frac{\partial}{\partial w} \sum_k (y_k - a_k)^2 = \sum_k \frac{\partial}{\partial w} (y_k - a_k)^2 \quad (18.10)$$

where the index  $k$  ranges over nodes in the output layer. Each term in the final summation is just the gradient of the loss for the  $k$ th output, computed as if the other outputs did not exist. Hence, we can decompose an  $m$ -output learning problem into  $m$  learning problems, provided we remember to add up the gradient contributions from each of them when updating the weights.

The major complication comes from the addition of hidden layers to the network. Whereas the error  $\mathbf{y} - \mathbf{h}_{\mathbf{w}}$  at the output layer is clear, the error at the hidden layers seems mysterious because the training data do not say what value the hidden nodes should have. Fortunately, it turns out that we can **back-propagate** the error from the output layer to the hidden layers. The back-propagation process emerges directly from a derivation of the overall error gradient. First, we will describe the process with an intuitive justification; then, we will show the derivation.

BACK-PROPAGATION

At the output layer, the weight-update rule is identical to Equation (18.8). We have multiple output units, so let  $Err_k$  be the  $k$ th component of the error vector  $\mathbf{y} - \mathbf{h}_{\mathbf{w}}$ . We will also find it convenient to define a modified error  $\Delta_k = Err_k \times g'(in_k)$ , so that the weight-update rule becomes

$$w_{j,k} \leftarrow w_{j,k} + \alpha \times a_j \times \Delta_k. \quad (18.11)$$

To update the connections between the input units and the hidden units, we need to define a quantity analogous to the error term for output nodes. Here is where we do the error back-propagation. The idea is that hidden node  $j$  is “responsible” for some fraction of the error  $\Delta_k$  in each of the output nodes to which it connects. Thus, the  $\Delta_k$  values are divided according to the strength of the connection between the hidden node and the output node and are propagated back to provide the  $\Delta_j$  values for the hidden layer. The propagation rule for the  $\Delta$  values is the following:

$$\Delta_j = g'(in_j) \sum_k w_{j,k} \Delta_k. \quad (18.12)$$

```

function BACK-PROP-LEARNING(examples, network) returns a neural network
  inputs: examples, a set of examples, each with input vector x and output vector y
           network, a multilayer network with  $L$  layers, weights  $w_{i,j}$ , activation function  $g$ 
  local variables:  $\Delta$ , a vector of errors, indexed by network node

  repeat
    for each weight  $w_{i,j}$  in network do
       $w_{i,j} \leftarrow$  a small random number
    for each example (x, y) in examples do
      /* Propagate the inputs forward to compute the outputs */
      for each node  $i$  in the input layer do
         $a_i \leftarrow x_i$ 
      for  $\ell = 2$  to  $L$  do
        for each node  $j$  in layer  $\ell$  do
           $in_j \leftarrow \sum_i w_{i,j} a_i$ 
           $a_j \leftarrow g(in_j)$ 
      /* Propagate deltas backward from output layer to input layer */
      for each node  $j$  in the output layer do
         $\Delta[j] \leftarrow g'(in_j) \times (y_j - a_j)$ 
      for  $\ell = L - 1$  to  $1$  do
        for each node  $i$  in layer  $\ell$  do
           $\Delta[i] \leftarrow g'(in_i) \sum_j w_{i,j} \Delta[j]$ 
      /* Update every weight in network using deltas */
      for each weight  $w_{i,j}$  in network do
         $w_{i,j} \leftarrow w_{i,j} + \alpha \times a_i \times \Delta[j]$ 
  until some stopping criterion is satisfied
  return network

```

**Figure 18.24** The back-propagation algorithm for learning in multilayer networks.

Now the weight-update rule for the weights between the inputs and the hidden layer is essentially identical to the update rule for the output layer:

$$w_{i,j} \leftarrow w_{i,j} + \alpha \times a_i \times \Delta_j .$$

The back-propagation process can be summarized as follows:

- Compute the  $\Delta$  values for the output units, using the observed error.
- Starting with output layer, repeat the following for each layer in the network, until the earliest hidden layer is reached:
  - Propagate the  $\Delta$  values back to the previous layer.
  - Update the weights between the two layers.

The detailed algorithm is shown in Figure 18.24.

For the mathematically inclined, we will now derive the back-propagation equations from first principles. The derivation is quite similar to the gradient calculation for logistic

regression (leading up to Equation (18.8) on page 727), except that we have to use the chain rule more than once.

Following Equation (18.10), we compute just the gradient for  $Loss_k = (y_k - a_k)^2$  at the  $k$ th output. The gradient of this loss with respect to weights connecting the hidden layer to the output layer will be zero except for weights  $w_{j,k}$  that connect to the  $k$ th output unit. For those weights, we have

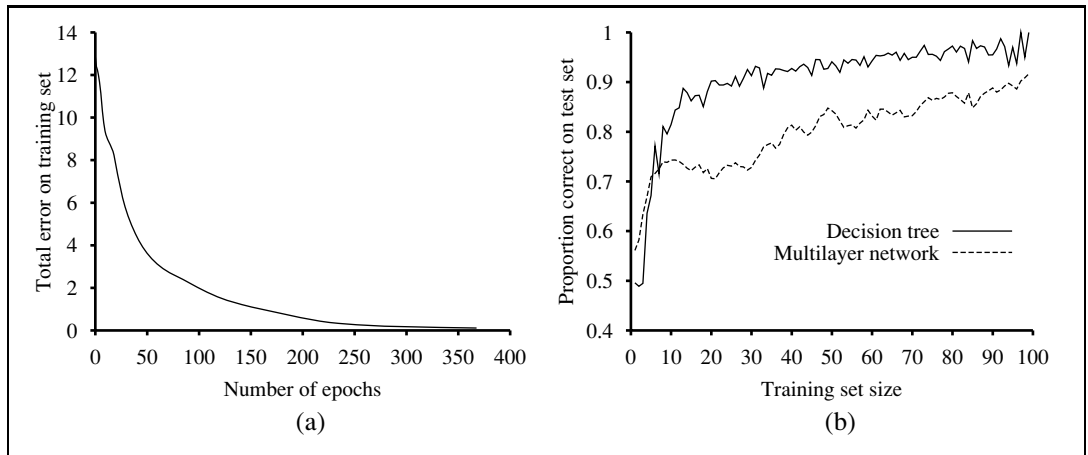
$$\begin{aligned} \frac{\partial Loss_k}{\partial w_{j,k}} &= -2(y_k - a_k) \frac{\partial a_k}{\partial w_{j,k}} = -2(y_k - a_k) \frac{\partial g(in_k)}{\partial w_{j,k}} \\ &= -2(y_k - a_k) g'(in_k) \frac{\partial in_k}{\partial w_{j,k}} = -2(y_k - a_k) g'(in_k) \frac{\partial}{\partial w_{j,k}} \left( \sum_j w_{j,k} a_j \right) \\ &= -2(y_k - a_k) g'(in_k) a_j = -a_j \Delta_k, \end{aligned}$$

with  $\Delta_k$  defined as before. To obtain the gradient with respect to the  $w_{i,j}$  weights connecting the input layer to the hidden layer, we have to expand out the activations  $a_j$  and reapply the chain rule. We will show the derivation in gory detail because it is interesting to see how the derivative operator propagates back through the network:

$$\begin{aligned} \frac{\partial Loss_k}{\partial w_{i,j}} &= -2(y_k - a_k) \frac{\partial a_k}{\partial w_{i,j}} = -2(y_k - a_k) \frac{\partial g(in_k)}{\partial w_{i,j}} \\ &= -2(y_k - a_k) g'(in_k) \frac{\partial in_k}{\partial w_{i,j}} = -2\Delta_k \frac{\partial}{\partial w_{i,j}} \left( \sum_j w_{j,k} a_j \right) \\ &= -2\Delta_k w_{j,k} \frac{\partial a_j}{\partial w_{i,j}} = -2\Delta_k w_{j,k} \frac{\partial g(in_j)}{\partial w_{i,j}} \\ &= -2\Delta_k w_{j,k} g'(in_j) \frac{\partial in_j}{\partial w_{i,j}} \\ &= -2\Delta_k w_{j,k} g'(in_j) \frac{\partial}{\partial w_{i,j}} \left( \sum_i w_{i,j} a_i \right) \\ &= -2\Delta_k w_{j,k} g'(in_j) a_i = -a_i \Delta_j, \end{aligned}$$

where  $\Delta_j$  is defined as before. Thus, we obtain the update rules obtained earlier from intuitive considerations. It is also clear that the process can be continued for networks with more than one hidden layer, which justifies the general algorithm given in Figure 18.24.

Having made it through (or skipped over) all the mathematics, let's see how a single-hidden-layer network performs on the restaurant problem. First, we need to determine the structure of the network. We have 10 attributes describing each example, so we will need 10 input units. Should we have one hidden layer or two? How many nodes in each layer? Should they be fully connected? There is no good theory that will tell us the answer. (See the next section.) As always, we can use cross-validation: try several different structures and see which one works best. It turns out that a network with one hidden layer containing four nodes is about right for this problem. In Figure 18.25, we show two curves. The first is a training curve showing the mean squared error on a given training set of 100 restaurant examples



**Figure 18.25** (a) Training curve showing the gradual reduction in error as weights are modified over several epochs, for a given set of examples in the restaurant domain. (b) Comparative learning curves showing that decision-tree learning does slightly better on the restaurant problem than back-propagation in a multilayer network.

during the weight-updating process. This demonstrates that the network does indeed converge to a perfect fit to the training data. The second curve is the standard learning curve for the restaurant data. The neural network does learn well, although not quite as fast as decision-tree learning; this is perhaps not surprising, because the data were generated from a simple decision tree in the first place.

Neural networks are capable of far more complex learning tasks of course, although it must be said that a certain amount of twiddling is needed to get the network structure right and to achieve convergence to something close to the global optimum in weight space. There are literally tens of thousands of published applications of neural networks. Section 18.11.1 looks at one such application in more depth.

### 18.7.5 Learning neural network structures

So far, we have considered the problem of learning weights, given a fixed network structure; just as with Bayesian networks, we also need to understand how to find the best network structure. If we choose a network that is too big, it will be able to memorize all the examples by forming a large lookup table, but will not necessarily generalize well to inputs that have not been seen before.<sup>10</sup> In other words, like all statistical models, neural networks are subject to **overfitting** when there are too many parameters in the model. We saw this in Figure 18.1 (page 696), where the high-parameter models in (b) and (c) fit all the data, but might not generalize as well as the low-parameter models in (a) and (d).

If we stick to fully connected networks, the only choices to be made concern the number

<sup>10</sup> It has been observed that very large networks *do* generalize well *as long as the weights are kept small*. This restriction keeps the activation values in the *linear* region of the sigmoid function  $g(x)$  where  $x$  is close to zero. This, in turn, means that the network behaves like a linear function (Exercise 18.22) with far fewer parameters.

of hidden layers and their sizes. The usual approach is to try several and keep the best. The **cross-validation** techniques of Chapter 18 are needed if we are to avoid **peeking** at the test set. That is, we choose the network architecture that gives the highest prediction accuracy on the validation sets.

OPTIMAL BRAIN  
DAMAGE

If we want to consider networks that are not fully connected, then we need to find some effective search method through the very large space of possible connection topologies. The **optimal brain damage** algorithm begins with a fully connected network and removes connections from it. After the network is trained for the first time, an information-theoretic approach identifies an optimal selection of connections that can be dropped. The network is then retrained, and if its performance has not decreased then the process is repeated. In addition to removing connections, it is also possible to remove units that are not contributing much to the result.

TILING

Several algorithms have been proposed for growing a larger network from a smaller one. One, the **tiling** algorithm, resembles decision-list learning. The idea is to start with a single unit that does its best to produce the correct output on as many of the training examples as possible. Subsequent units are added to take care of the examples that the first unit got wrong. The algorithm adds only as many units as are needed to cover all the examples.

## 18.8 NONPARAMETRIC MODELS

Linear regression and neural networks use the training data to estimate a fixed set of parameters  $\mathbf{w}$ . That defines our hypothesis  $h_{\mathbf{w}}(\mathbf{x})$ , and at that point we can throw away the training data, because they are all summarized by  $\mathbf{w}$ . A learning model that summarizes data with a set of parameters of fixed size (independent of the number of training examples) is called a **parametric model**.

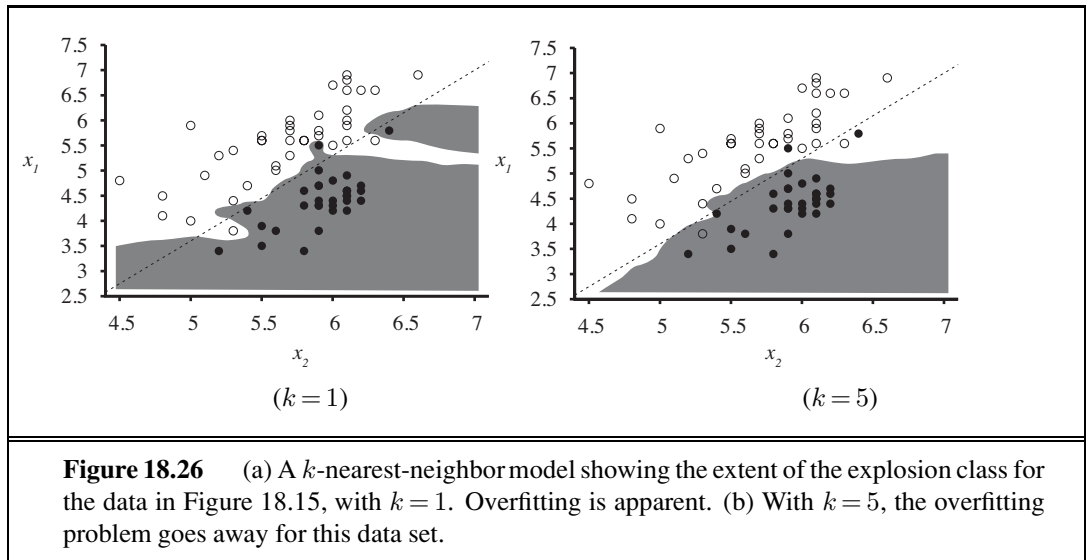
PARAMETRIC MODEL

No matter how much data you throw at a parametric model, it won't change its mind about how many parameters it needs. When data sets are small, it makes sense to have a strong restriction on the allowable hypotheses, to avoid overfitting. But when there are thousands or millions or billions of examples to learn from, it seems like a better idea to let the data speak for themselves rather than forcing them to speak through a tiny vector of parameters. If the data say that the correct answer is a very wiggly function, we shouldn't restrict ourselves to linear or slightly wiggly functions.

NONPARAMETRIC  
MODEL

A **nonparametric model** is one that cannot be characterized by a bounded set of parameters. For example, suppose that each hypothesis we generate simply retains within itself all of the training examples and uses all of them to predict the next example. Such a hypothesis family would be nonparametric because the effective number of parameters is unbounded—it grows with the number of examples. This approach is called **instance-based learning** or **memory-based learning**. The simplest instance-based learning method is **table lookup**: take all the training examples, put them in a lookup table, and then when asked for  $h(\mathbf{x})$ , see if  $\mathbf{x}$  is in the table; if it is, return the corresponding  $y$ . The problem with this method is that it does not generalize well: when  $\mathbf{x}$  is not in the table all it can do is return some default value.

INSTANCE-BASED  
LEARNING  
TABLE LOOKUP



### 18.8.1 Nearest neighbor models

NEAREST  
NEIGHBORS

We can improve on table lookup with a slight variation: given a query  $\mathbf{x}_q$ , find the  $k$  examples that are *nearest* to  $\mathbf{x}_q$ . This is called  **$k$ -nearest neighbors** lookup. We'll use the notation  $NN(k, \mathbf{x}_q)$  to denote the set of  $k$  nearest neighbors.

To do classification, first find  $NN(k, \mathbf{x}_q)$ , then take the plurality vote of the neighbors (which is the majority vote in the case of binary classification). To avoid ties,  $k$  is always chosen to be an odd number. To do regression, we can take the mean or median of the  $k$  neighbors, or we can solve a linear regression problem on the neighbors.

In Figure 18.26, we show the decision boundary of  $k$ -nearest-neighbors classification for  $k = 1$  and 5 on the earthquake data set from Figure 18.15. Nonparametric methods are still subject to underfitting and overfitting, just like parametric methods. In this case 1-nearest neighbors is overfitting; it reacts too much to the black outlier in the upper right and the white outlier at (5.4, 3.7). The 5-nearest-neighbors decision boundary is good; higher  $k$  would underfit. As usual, cross-validation can be used to select the best value of  $k$ .

MINKOWSKI  
DISTANCE

The very word “nearest” implies a distance metric. How do we measure the distance from a query point  $\mathbf{x}_q$  to an example point  $\mathbf{x}_j$ ? Typically, distances are measured with a **Minkowski distance** or  $L^p$  norm, defined as

$$L^p(\mathbf{x}_j, \mathbf{x}_q) = \left( \sum_i |x_{j,i} - x_{q,i}|^p \right)^{1/p}.$$

HAMMING DISTANCE

With  $p = 2$  this is Euclidean distance and with  $p = 1$  it is Manhattan distance. With Boolean attribute values, the number of attributes on which the two points differ is called the **Hamming distance**. Often  $p = 2$  is used if the dimensions are measuring similar properties, such as the width, height and depth of parts on a conveyor belt, and Manhattan distance is used if they are dissimilar, such as age, weight, and gender of a patient. Note that if we use the raw numbers from each dimension then the total distance will be affected by a change in scale in any dimension. That is, if we change dimension  $i$  from measurements in centimeters to

NORMALIZATION

miles while keeping the other dimensions the same, we'll get different nearest neighbors. To avoid this, it is common to apply **normalization** to the measurements in each dimension. One simple approach is to compute the mean  $\mu_i$  and standard deviation  $\sigma_i$  of the values in each dimension, and rescale them so that  $x_{j,i}$  becomes  $(x_{j,i} - \mu_i)/\sigma_i$ . A more complex metric known as the **Mahalanobis distance** takes into account the covariance between dimensions.

MAHALANOBIS  
DISTANCE

In low-dimensional spaces with plenty of data, nearest neighbors works very well: we are likely to have enough nearby data points to get a good answer. But as the number of dimensions rises we encounter a problem: the nearest neighbors in high-dimensional spaces are usually not very near! Consider  $k$ -nearest-neighbors on a data set of  $N$  points uniformly distributed throughout the interior of an  $n$ -dimensional unit hypercube. We'll define the  $k$ -neighborhood of a point as the smallest hypercube that contains the  $k$ -nearest neighbors. Let  $\ell$  be the average side length of a neighborhood. Then the volume of the neighborhood (which contains  $k$  points) is  $\ell^n$  and the volume of the full cube (which contains  $N$  points) is 1. So, on average,  $\ell^n = k/N$ . Taking  $n$ th roots of both sides we get  $\ell = (k/N)^{1/n}$ .

CURSE OF  
DIMENSIONALITY

To be concrete, let  $k = 10$  and  $N = 1,000,000$ . In two dimensions ( $n = 2$ ; a unit square), the average neighborhood has  $\ell = 0.003$ , a small fraction of the unit square, and in 3 dimensions  $\ell$  is just 2% of the edge length of the unit cube. But by the time we get to 17 dimensions,  $\ell$  is half the edge length of the unit hypercube, and in 200 dimensions it is 94%. This problem has been called the **curse of dimensionality**.

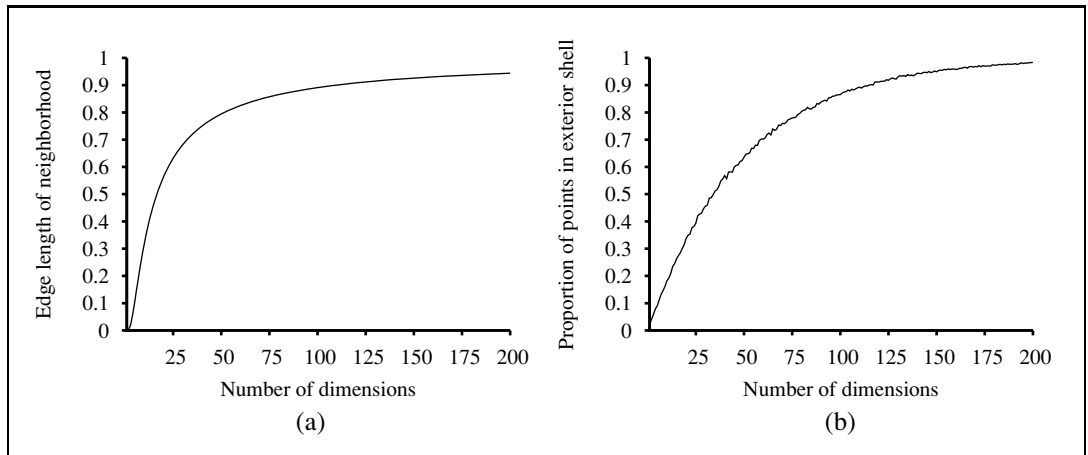
Another way to look at it: consider the points that fall within a thin shell making up the outer 1% of the unit hypercube. These are outliers; in general it will be hard to find a good value for them because we will be extrapolating rather than interpolating. In one dimension, these outliers are only 2% of the points on the unit line (those points where  $x < .01$  or  $x > .99$ ), but in 200 dimensions, over 98% of the points fall within this thin shell—almost all the points are outliers. You can see an example of a poor nearest-neighbors fit on outliers if you look ahead to Figure 18.28(b).

The  $NN(k, \mathbf{x}_q)$  function is conceptually trivial: given a set of  $N$  examples and a query  $\mathbf{x}_q$ , iterate through the examples, measure the distance to  $\mathbf{x}_q$  from each one, and keep the best  $k$ . If we are satisfied with an implementation that takes  $O(N)$  execution time, then that is the end of the story. But instance-based methods are designed for large data sets, so we would like an algorithm with sublinear run time. Elementary analysis of algorithms tells us that exact table lookup is  $O(N)$  with a sequential table,  $O(\log N)$  with a binary tree, and  $O(1)$  with a hash table. We will now see that binary trees and hash tables are also applicable for finding nearest neighbors.

## 18.8.2 Finding nearest neighbors with k-d trees

K-D TREE

A balanced binary tree over data with an arbitrary number of dimensions is called a **k-d tree**, for k-dimensional tree. (In our notation, the number of dimensions is  $n$ , so they would be  $n$ -d trees. The construction of a k-d tree is similar to the construction of a one-dimensional balanced binary tree. We start with a set of examples and at the root node we split them along the  $i$ th dimension by testing whether  $x_i \leq m$ . We chose the value  $m$  to be the median of the examples along the  $i$ th dimension; thus half the examples will be in the left branch of the tree



**Figure 18.27** The curse of dimensionality: (a) The length of the average neighborhood for 10-nearest-neighbors in a unit hypercube with 1,000,000 points, as a function of the number of dimensions. (b) The proportion of points that fall within a thin shell consisting of the outer 1% of the hypercube, as a function of the number of dimensions. Sampled from 10,000 randomly distributed points.

and half in the right. We then recursively make a tree for the left and right sets of examples, stopping when there are fewer than two examples left. To choose a dimension to split on at each node of the tree, one can simply select dimension  $i \bmod n$  at level  $i$  of the tree. (Note that we may need to split on any given dimension several times as we proceed down the tree.) Another strategy is to split on the dimension that has the widest spread of values.

Exact lookup from a k-d tree is just like lookup from a binary tree (with the slight complication that you need to pay attention to which dimension you are testing at each node). But nearest neighbor lookup is more complicated. As we go down the branches, splitting the examples in half, in some cases we can discard the other half of the examples. But not always. Sometimes the point we are querying for falls very close to the dividing boundary. The query point itself might be on the left hand side of the boundary, but one or more of the  $k$  nearest neighbors might actually be on the right-hand side. We have to test for this possibility by computing the distance of the query point to the dividing boundary, and then searching both sides if we can't find  $k$  examples on the left that are closer than this distance. Because of this problem, k-d trees are appropriate only when there are many more examples than dimensions, preferably at least  $2^n$  examples. Thus, k-d trees work well with up to 10 dimensions with thousands of examples or up to 20 dimensions with millions of examples. If we don't have enough examples, lookup is no faster than a linear scan of the entire data set.

### 18.8.3 Locality-sensitive hashing

Hash tables have the potential to provide even faster lookup than binary trees. But how can we find nearest neighbors using a hash table, when hash codes rely on an *exact* match? Hash codes randomly distribute values among the bins, but we want to have near points grouped together in the same bin; we want a **locality-sensitive hash** (LSH).



We can't use hashes to solve  $NN(k, \mathbf{x}_q)$  exactly, but with a clever use of randomized algorithms, we can find an *approximate* solution. First we define the **approximate near-neighbors** problem: given a data set of example points and a query point  $\mathbf{x}_q$ , find, with high probability, an example point (or points) that is near  $\mathbf{x}_q$ . To be more precise, we require that if there is a point  $\mathbf{x}_j$  that is within a radius  $r$  of  $\mathbf{x}_q$ , then with high probability the algorithm will find a point  $\mathbf{x}_{j'}$  that is within distance  $cr$  of  $q$ . If there is no point within radius  $r$  then the algorithm is allowed to report failure. The values of  $c$  and “high probability” are parameters of the algorithm.

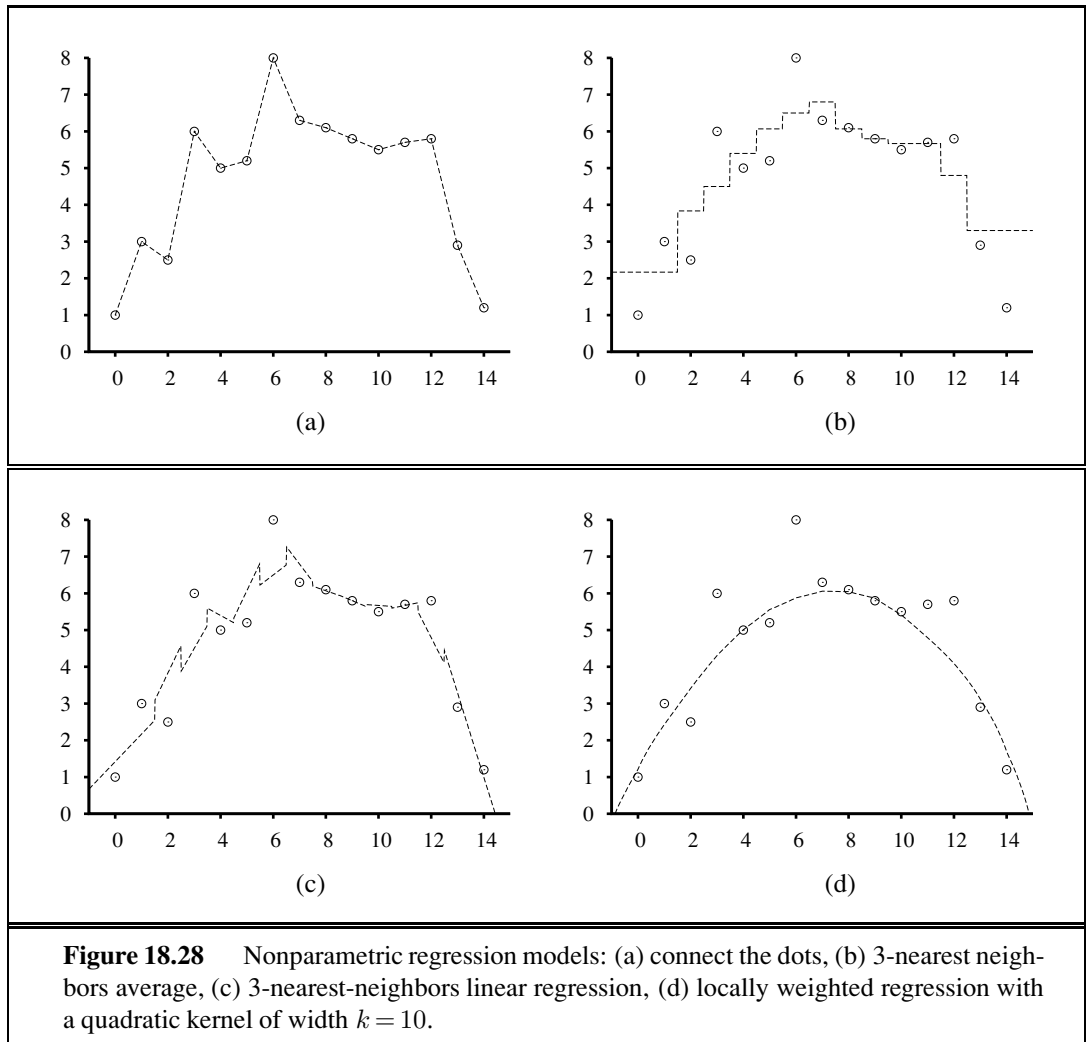
To solve approximate near neighbors, we will need a hash function  $g(\mathbf{x})$  that has the property that, for any two points  $\mathbf{x}_j$  and  $\mathbf{x}_{j'}$ , the probability that they have the same hash code is small if their distance is more than  $cr$ , and is high if their distance is less than  $r$ . For simplicity we will treat each point as a bit string. (Any features that are not Boolean can be encoded into a set of Boolean features.)

The intuition we rely on is that if two points are close together in an  $n$ -dimensional space, then they will necessarily be close when projected down onto a one-dimensional space (a line). In fact, we can discretize the line into bins—hash buckets—so that, with high probability, near points project down to exactly the same bin. Points that are far away from each other will tend to project down into different bins for most projections, but there will always be a few projections that coincidentally project far-apart points into the same bin. Thus, the bin for point  $\mathbf{x}_q$  contains many (but not all) points that are near to  $\mathbf{x}_q$ , as well as some points that are far away.

The trick of LSH is to create *multiple* random projections and combine them. A random projection is just a random subset of the bit-string representation. We choose  $\ell$  different random projections and create  $\ell$  hash tables,  $g_1(\mathbf{x}), \dots, g_\ell(\mathbf{x})$ . We then enter all the examples into each hash table. Then when given a query point  $\mathbf{x}_q$ , we fetch the set of points in bin  $g_k(q)$  for each  $k$ , and union these sets together into a set of candidate points,  $C$ . Then we compute the actual distance to  $\mathbf{x}_q$  for each of the points in  $C$  and return the  $k$  closest points. With high probability, each of the points that are near to  $\mathbf{x}_q$  will show up in at least one of the bins, and although some far-away points will show up as well, we can ignore those. With large real-world problems, such as finding the near neighbors in a data set of 13 million Web images using 512 dimensions (Torralba *et al.*, 2008), locality-sensitive hashing needs to examine only a few thousand images out of 13 million to find nearest neighbors; a thousand-fold speedup over exhaustive or k-d tree approaches.

#### 18.8.4 Nonparametric regression

Now we'll look at nonparametric approaches to *regression* rather than classification. Figure 18.28 shows an example of some different models. In (a), we have perhaps the simplest method of all, known informally as “connect-the-dots,” and superciliously as “piecewise-linear nonparametric regression.” This model creates a function  $h(x)$  that, when given a query  $x_q$ , solves the ordinary linear regression problem with just two points: the training examples immediately to the left and right of  $x_q$ . When noise is low, this trivial method is actually not too bad, which is why it is a standard feature of charting software in spreadsheets.



**Figure 18.28** Nonparametric regression models: (a) connect the dots, (b) 3-nearest neighbors average, (c) 3-nearest-neighbor linear regression, (d) locally weighted regression with a quadratic kernel of width  $k=10$ .

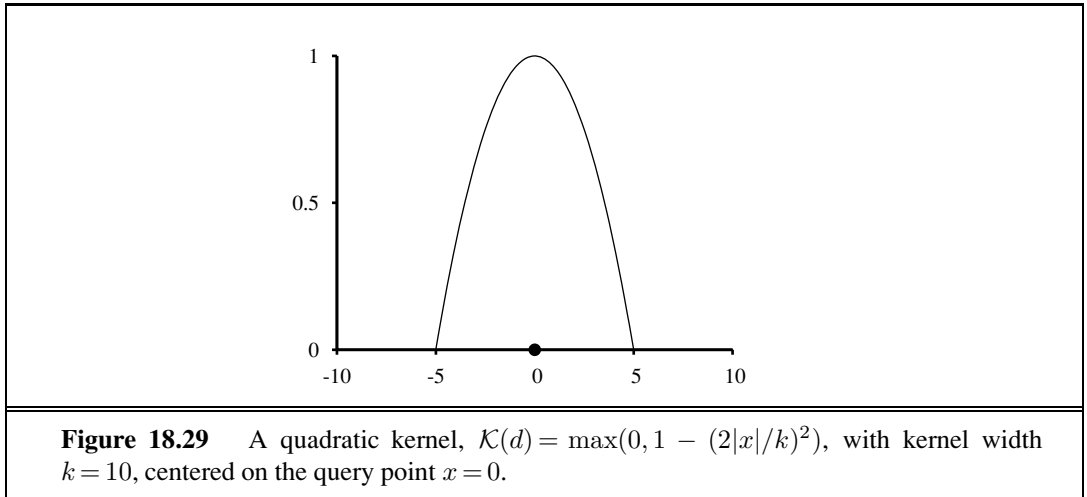
But when the data are noisy, the resulting function is spiky, and does not generalize well.

**$k$ -nearest-neighbors regression** (Figure 18.28(b)) improves on connect-the-dots. Instead of using just the two examples to the left and right of a query point  $x_q$ , we use the  $k$  nearest neighbors (here 3). A larger value of  $k$  tends to smooth out the magnitude of the spikes, although the resulting function has discontinuities. In (b), we have the  $k$ -nearest-neighbors average:  $h(x)$  is the mean value of the  $k$  points,  $\sum y_j/k$ . Notice that at the outlying points, near  $x=0$  and  $x=14$ , the estimates are poor because all the evidence comes from one side (the interior), and ignores the trend. In (c), we have  $k$ -nearest-neighbor linear regression, which finds the best line through the  $k$  examples. This does a better job of capturing trends at the outliers, but is still discontinuous. In both (b) and (c), we're left with the question of how to choose a good value for  $k$ . The answer, as usual, is cross-validation.

**Locally weighted regression** (Figure 18.28(d)) gives us the advantages of nearest neighbors, without the discontinuities. To avoid discontinuities in  $h(x)$ , we need to avoid disconti-

NEAREST-  
NEIGHBORS  
REGRESSION

LOCALLY WEIGHTED  
REGRESSION



nities in the set of examples we use to estimate  $h(x)$ . The idea of locally weighted regression is that at each query point  $x_q$ , the examples that are close to  $x_q$  are weighted heavily, and the examples that are farther away are weighted less heavily or not at all. The decrease in weight over distance is always gradual, not sudden.

KERNEL

We decide how much to weight each example with a function known as a **kernel**. A kernel function looks like a bump; in Figure 18.29 we see the specific kernel used to generate Figure 18.28(d). We can see that the weight provided by this kernel is highest in the center and reaches zero at a distance of  $\pm 5$ . Can we choose just any function for a kernel? No. First, note that we invoke a kernel function  $\mathcal{K}$  with  $\mathcal{K}(\text{Distance}(\mathbf{x}_j, \mathbf{x}_q))$ , where  $\mathbf{x}_q$  is a query point that is a given distance from  $\mathbf{x}_j$ , and we want to know how much to weight that distance. So  $\mathcal{K}$  should be symmetric around 0 and have a maximum at 0. The area under the kernel must remain bounded as we go to  $\pm\infty$ . Other shapes, such as Gaussians, have been used for kernels, but the latest research suggests that the choice of shape doesn't matter much. We do have to be careful about the width of the kernel. Again, this is a parameter of the model that is best chosen by cross-validation. Just as in choosing the  $k$  for nearest neighbors, if the kernels are too wide we'll get underfitting and if they are too narrow we'll get overfitting. In Figure 18.29(d), the value of  $k = 10$  gives a smooth curve that looks about right—but maybe it does not pay enough attention to the outlier at  $x = 6$ ; a narrower kernel width would be more responsive to individual points.

Doing locally weighted regression with kernels is now straightforward. For a given query point  $\mathbf{x}_q$  we solve the following weighted regression problem using gradient descent:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_j \mathcal{K}(\text{Distance}(\mathbf{x}_q, \mathbf{x}_j)) (y_j - \mathbf{w} \cdot \mathbf{x}_j)^2,$$

where  $\text{Distance}$  is any of the distance metrics discussed for nearest neighbors. Then the answer is  $h(\mathbf{x}_q) = \mathbf{w}^* \cdot \mathbf{x}_q$ .

Note that we need to solve a new regression problem for *every* query point—that's what it means to be *local*. (In ordinary linear regression, we solved the regression problem once, globally, and then used the same  $h_{\mathbf{w}}$  for any query point.) Mitigating against this extra work

is the fact that each regression problem will be easier to solve, because it involves only the examples with nonzero weight—the examples whose kernels overlap the query point. When kernel widths are small, this may be just a few points.

Most nonparametric models have the advantage that it is easy to do leave-one-out cross-validation without having to recompute everything. With a  $k$ -nearest-neighbors model, for instance, when given a test example  $(\mathbf{x}, y)$  we retrieve the  $k$  nearest neighbors once, compute the per-example loss  $L(y, h(\mathbf{x}))$  from them, and record that as the leave-one-out result for every example that is not one of the neighbors. Then we retrieve the  $k + 1$  nearest neighbors and record distinct results for leaving out each of the  $k$  neighbors. With  $N$  examples the whole process is  $O(k)$ , not  $O(kN)$ .

## 18.9 SUPPORT VECTOR MACHINES

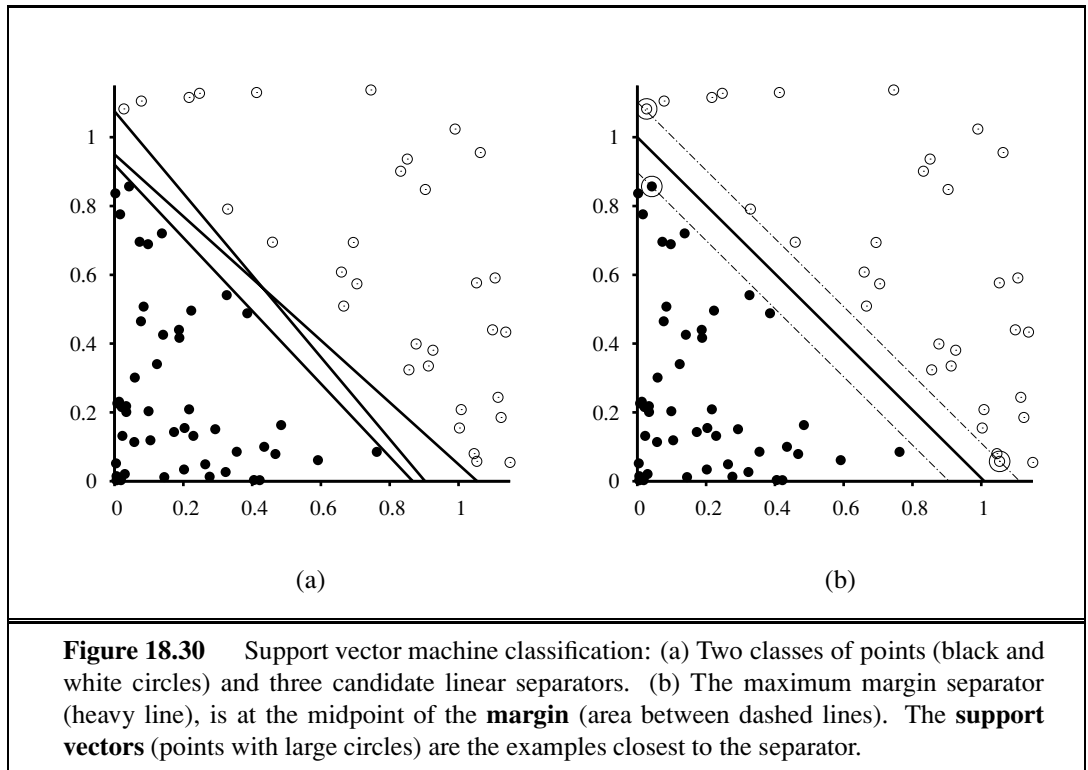
### SUPPORT VECTOR MACHINE

The **support vector machine** or SVM framework is currently the most popular approach for “off-the-shelf” supervised learning: if you don’t have any specialized prior knowledge about a domain, then the SVM is an excellent method to try first. There are three properties that make SVMs attractive:

1. SVMs construct a **maximum margin separator**—a decision boundary with the largest possible distance to example points. This helps them generalize well.
2. SVMs create a linear separating hyperplane, but they have the ability to embed the data into a higher-dimensional space, using the so-called **kernel trick**. Often, data that are not linearly separable in the original input space are easily separable in the higher-dimensional space. The high-dimensional linear separator is actually nonlinear in the original space. This means the hypothesis space is greatly expanded over methods that use strictly linear representations.
3. SVMs are a nonparametric method—they retain training examples and potentially need to store them all. On the other hand, in practice they often end up retaining only a small fraction of the number of examples—sometimes as few as a small constant times the number of dimensions. Thus SVMs combine the advantages of nonparametric and parametric models: they have the flexibility to represent complex functions, but they are resistant to overfitting.

You could say that SVMs are successful because of one key insight and one neat trick. We will cover each in turn. In Figure 18.30(a), we have a binary classification problem with three candidate decision boundaries, each a linear separator. Each of them is consistent with all the examples, so from the point of view of 0/1 loss, each would be equally good. Logistic regression would find some separating line; the exact location of the line depends on *all* the example points. The key insight of SVMs is that some examples are more important than others, and that paying attention to them can lead to better generalization.

Consider the lowest of the three separating lines in (a). It comes very close to 5 of the black examples. Although it classifies all the examples correctly, and thus minimizes loss, it



**Figure 18.30** Support vector machine classification: (a) Two classes of points (black and white circles) and three candidate linear separators. (b) The maximum margin separator (heavy line), is at the midpoint of the **margin** (area between dashed lines). The **support vectors** (points with large circles) are the examples closest to the separator.

should make you nervous that so many examples are close to the line; it may be that other black examples will turn out to fall on the other side of the line.

SVMs address this issue: Instead of minimizing expected *empirical loss* on the training data, SVMs attempt to minimize expected *generalization loss*. We don't know where the as-yet-unseen points may fall, but under the probabilistic assumption that they are drawn from the same distribution as the previously seen examples, there are some arguments from computational learning theory (Section 18.5) suggesting that we minimize generalization loss by choosing the separator that is farthest away from the examples we have seen so far. We call this separator, shown in Figure 18.30(b) the **maximum margin separator**. The **margin** is the width of the area bounded by dashed lines in the figure—twice the distance from the separator to the nearest example point.

Now, how do we find this separator? Before showing the equations, some notation: Traditionally SVMs use the convention that class labels are +1 and -1, instead of the +1 and 0 we have been using so far. Also, where we put the intercept into the weight vector  $\mathbf{w}$  (and a corresponding dummy 1 value into  $x_{j,0}$ ), SVMs do not do that; they keep the intercept as a separate parameter,  $b$ . With that in mind, the separator is defined as the set of points  $\{\mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = 0\}$ . We could search the space of  $\mathbf{w}$  and  $b$  with gradient descent to find the parameters that maximize the margin while correctly classifying all the examples.

However, it turns out there is another approach to solving this problem. We won't show the details, but will just say that there is an alternative representation called the dual

MAXIMUM MARGIN  
SEPARATOR  
MARGIN

representation, in which the optimal solution is found by solving

$$\operatorname{argmax}_{\alpha} \sum_j \alpha_j - \frac{1}{2} \sum_{j,k} \alpha_j \alpha_k y_j y_k (\mathbf{x}_j \cdot \mathbf{x}_k) \quad (18.13)$$

QUADRATIC  
PROGRAMMING



subject to the constraints  $\alpha_j \geq 0$  and  $\sum_j \alpha_j y_j = 0$ . This is a **quadratic programming** optimization problem, for which there are good software packages. Once we have found the vector  $\alpha$  we can get back to  $\mathbf{w}$  with the equation  $\mathbf{w} = \sum_j \alpha_j \mathbf{x}_j$ , or we can stay in the dual representation. There are three important properties of Equation (18.13). First, the expression is convex; it has a single global maximum that can be found efficiently. Second, *the data enter the expression only in the form of dot products of pairs of points*. This second property is also true of the equation for the separator itself; once the optimal  $\alpha_j$  have been calculated, it is

$$h(\mathbf{x}) = \operatorname{sign} \left( \sum_j \alpha_j y_j (\mathbf{x} \cdot \mathbf{x}_j) - b \right). \quad (18.14)$$

SUPPORT VECTOR

A final important property is that the weights  $\alpha_j$  associated with each data point are *zero* except for the **support vectors**—the points closest to the separator. (They are called “support” vectors because they “hold up” the separating plane.) Because there are usually many fewer support vectors than examples, SVMs gain some of the advantages of parametric models.

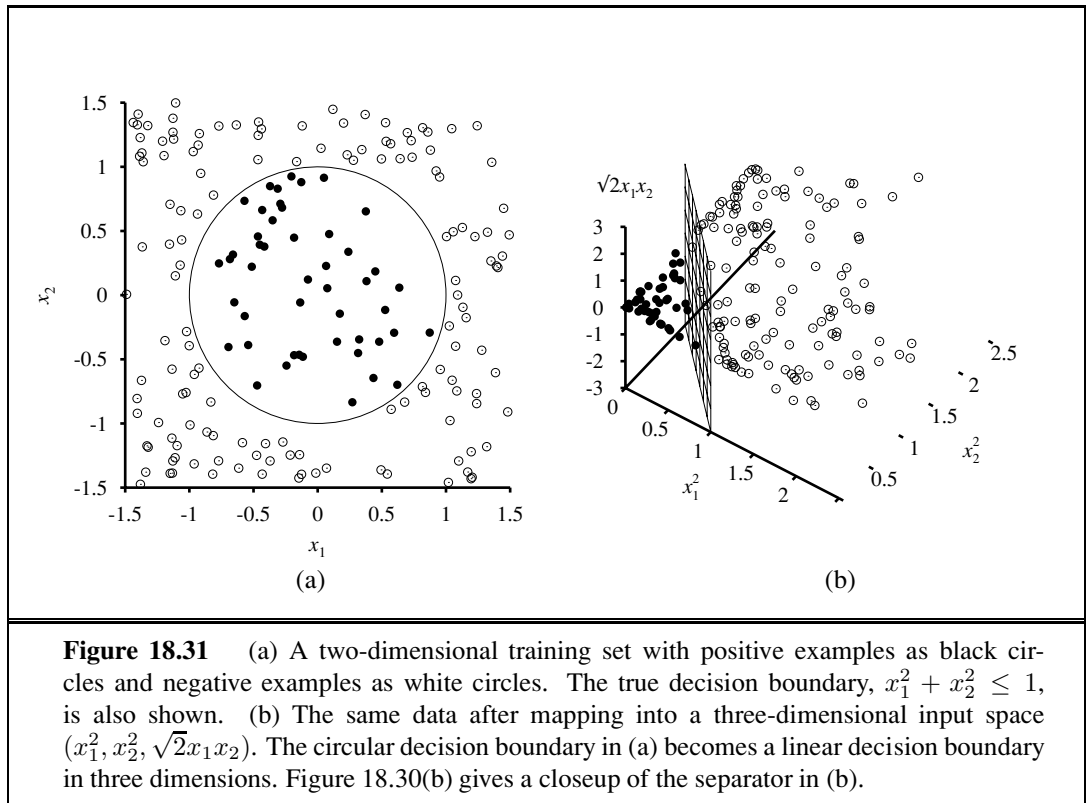
What if the examples are not linearly separable? Figure 18.31(a) shows an input space defined by attributes  $\mathbf{x} = (x_1, x_2)$ , with positive examples ( $y = +1$ ) inside a circular region and negative examples ( $y = -1$ ) outside. Clearly, there is no linear separator for this problem. Now, suppose we re-express the input data—i.e., we map each input vector  $\mathbf{x}$  to a new vector of feature values,  $F(\mathbf{x})$ . In particular, let us use the three features

$$f_1 = x_1^2, \quad f_2 = x_2^2, \quad f_3 = \sqrt{2}x_1x_2. \quad (18.15)$$

We will see shortly where these came from, but for now, just look at what happens. Figure 18.31(b) shows the data in the new, three-dimensional space defined by the three features; the data are *linearly separable* in this space! This phenomenon is actually fairly general: if data are mapped into a space of sufficiently high dimension, then they will almost always be linearly separable—if you look at a set of points from enough directions, you’ll find a way to make them line up. Here, we used only three dimensions;<sup>11</sup> Exercise 18.16 asks you to show that four dimensions suffice for linearly separating a circle anywhere in the plane (not just at the origin), and five dimensions suffice to linearly separate any ellipse. In general (with some special cases excepted) if we have  $N$  data points then they will always be separable in spaces of  $N - 1$  dimensions or more (Exercise 18.25).

Now, we would not usually expect to find a linear separator in the input space  $\mathbf{x}$ , but we can find linear separators in the high-dimensional feature space  $F(\mathbf{x})$  simply by replacing  $\mathbf{x}_j \cdot \mathbf{x}_k$  in Equation (18.13) with  $F(\mathbf{x}_j) \cdot F(\mathbf{x}_k)$ . This by itself is not remarkable—replacing  $\mathbf{x}$  by  $F(\mathbf{x})$  in *any* learning algorithm has the required effect—but the dot product has some special properties. It turns out that  $F(\mathbf{x}_j) \cdot F(\mathbf{x}_k)$  can often be computed without first computing  $F$

<sup>11</sup> The reader may notice that we could have used just  $f_1$  and  $f_2$ , but the 3D mapping illustrates the idea better.



**Figure 18.31** (a) A two-dimensional training set with positive examples as black circles and negative examples as white circles. The true decision boundary,  $x_1^2 + x_2^2 \leq 1$ , is also shown. (b) The same data after mapping into a three-dimensional input space  $(x_1^2, x_2^2, \sqrt{2}x_1x_2)$ . The circular decision boundary in (a) becomes a linear decision boundary in three dimensions. Figure 18.30(b) gives a closeup of the separator in (b).

for each point. In our three-dimensional feature space defined by Equation (18.15), a little bit of algebra shows that

$$F(\mathbf{x}_j) \cdot F(\mathbf{x}_k) = (\mathbf{x}_j \cdot \mathbf{x}_k)^2.$$

KERNEL FUNCTION

(That's why the  $\sqrt{2}$  is in  $f_3$ .) The expression  $(\mathbf{x}_j \cdot \mathbf{x}_k)^2$  is called a **kernel function**,<sup>12</sup> and is usually written as  $K(\mathbf{x}_j, \mathbf{x}_k)$ . The kernel function can be applied to pairs of input data to evaluate dot products in some corresponding feature space. So, we can find linear separators in the higher-dimensional feature space  $F(\mathbf{x})$  simply by replacing  $\mathbf{x}_j \cdot \mathbf{x}_k$  in Equation (18.13) with a kernel function  $K(\mathbf{x}_j, \mathbf{x}_k)$ . Thus, we can learn in the higher-dimensional space, but we compute only kernel functions rather than the full list of features for each data point.

MERCER'S THEOREM

POLYNOMIAL  
KERNEL

The next step is to see that there's nothing special about the kernel  $K(\mathbf{x}_j, \mathbf{x}_k) = (\mathbf{x}_j \cdot \mathbf{x}_k)^2$ . It corresponds to a particular higher-dimensional feature space, but other kernel functions correspond to other feature spaces. A venerable result in mathematics, **Mercer's theorem** (1909), tells us that any "reasonable"<sup>13</sup> kernel function corresponds to *some* feature space. These feature spaces can be very large, even for innocuous-looking kernels. For example, the **polynomial kernel**,  $K(\mathbf{x}_j, \mathbf{x}_k) = (1 + \mathbf{x}_j \cdot \mathbf{x}_k)^d$ , corresponds to a feature space whose dimension is exponential in  $d$ .

<sup>12</sup> This usage of "kernel function" is slightly different from the kernels in locally weighted regression. Some SVM kernels are distance metrics, but not all are.

<sup>13</sup> Here, "reasonable" means that the matrix  $\mathbf{K}_{jk} = K(\mathbf{x}_j, \mathbf{x}_k)$  is positive definite.

## KERNEL TRICK



This then is the clever **kernel trick**: Plugging these kernels into Equation (18.13), *optimal linear separators can be found efficiently in feature spaces with billions of (or, in some cases, infinitely many) dimensions*. The resulting linear separators, when mapped back to the original input space, can correspond to arbitrarily wiggly, nonlinear decision boundaries between the positive and negative examples.

## SOFT MARGIN

In the case of inherently noisy data, we may not want a linear separator in some high-dimensional space. Rather, we'd like a decision surface in a lower-dimensional space that does not cleanly separate the classes, but reflects the reality of the noisy data. That is possible with the **soft margin** classifier, which allows examples to fall on the wrong side of the decision boundary, but assigns them a penalty proportional to the distance required to move them back on the correct side.

## KERNELIZATION

The kernel method can be applied not only with learning algorithms that find optimal linear separators, but also with any other algorithm that can be reformulated to work only with dot products of pairs of data points, as in Equations 18.13 and 18.14. Once this is done, the dot product is replaced by a kernel function and we have a **kernelized** version of the algorithm. This can be done easily for  $k$ -nearest-neighbors and perceptron learning (Section 18.7.2), among others.

## 18.10 ENSEMBLE LEARNING

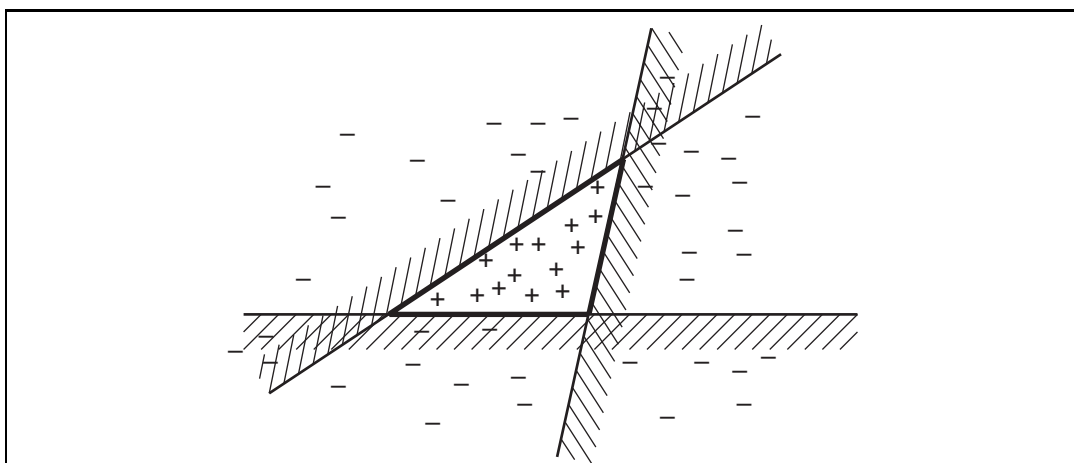
ENSEMBLE  
LEARNING

So far we have looked at learning methods in which a single hypothesis, chosen from a hypothesis space, is used to make predictions. The idea of **ensemble learning** methods is to select a collection, or **ensemble**, of hypotheses from the hypothesis space and combine their predictions. For example, during cross-validation we might generate twenty different decision trees, and have them vote on the best classification for a new example.

The motivation for ensemble learning is simple. Consider an ensemble of  $K = 5$  hypotheses and suppose that we combine their predictions using simple majority voting. For the ensemble to misclassify a new example, *at least three of the five hypotheses have to misclassify it*. The hope is that this is much less likely than a misclassification by a single hypothesis. Suppose we assume that each hypothesis  $h_k$  in the ensemble has an error of  $p$ —that is, the probability that a randomly chosen example is misclassified by  $h_k$  is  $p$ . Furthermore, suppose we assume that the errors made by each hypothesis are *independent*. In that case, if  $p$  is small, then the probability of a large number of misclassifications occurring is minuscule. For example, a simple calculation (Exercise 18.18) shows that using an ensemble of five hypotheses reduces an error rate of 1 in 10 down to an error rate of less than 1 in 100. Now, obviously the assumption of independence is unreasonable, because hypotheses are likely to be misled in the same way by any misleading aspects of the training data. But if the hypotheses are at least a little bit different, thereby reducing the correlation between their errors, then ensemble learning can be very useful.

Another way to think about the ensemble idea is as a generic way of enlarging the hypothesis space. That is, think of the ensemble itself as a hypothesis and the new hypothesis





**Figure 18.32** Illustration of the increased expressive power obtained by ensemble learning. We take three linear threshold hypotheses, each of which classifies positively on the unshaded side, and classify as positive any example classified positively by all three. The resulting triangular region is a hypothesis not expressible in the original hypothesis space.

space as the set of all possible ensembles constructable from hypotheses in the original space. Figure 18.32 shows how this can result in a more expressive hypothesis space. If the original hypothesis space allows for a simple and efficient learning algorithm, then the ensemble method provides a way to learn a much more expressive class of hypotheses without incurring much additional computational or algorithmic complexity.

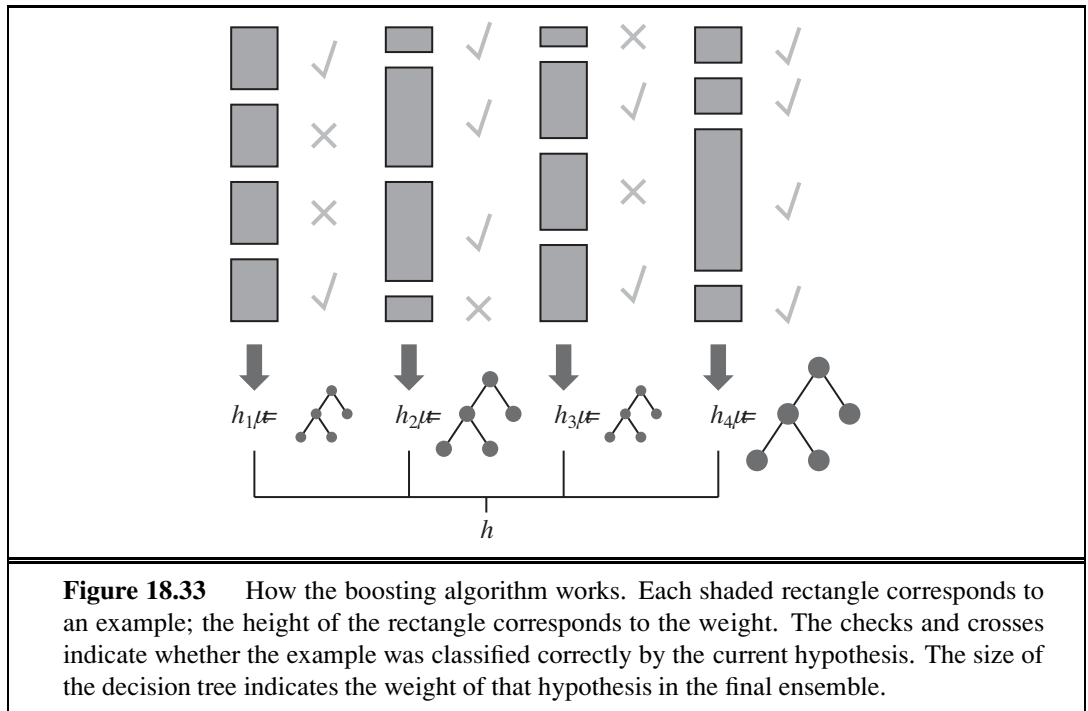
BOOSTING  
WEIGHTED TRAINING  
SET

The most widely used ensemble method is called **boosting**. To understand how it works, we need first to explain the idea of a **weighted training set**. In such a training set, each example has an associated weight  $w_j \geq 0$ . The higher the weight of an example, the higher is the importance attached to it during the learning of a hypothesis. It is straightforward to modify the learning algorithms we have seen so far to operate with weighted training sets.<sup>14</sup>

Boosting starts with  $w_j = 1$  for all the examples (i.e., a normal training set). From this set, it generates the first hypothesis,  $h_1$ . This hypothesis will classify some of the training examples correctly and some incorrectly. We would like the next hypothesis to do better on the misclassified examples, so we increase their weights while decreasing the weights of the correctly classified examples. From this new weighted training set, we generate hypothesis  $h_2$ . The process continues in this way until we have generated  $K$  hypotheses, where  $K$  is an input to the boosting algorithm. The final ensemble hypothesis is a weighted-majority combination of all the  $K$  hypotheses, each weighted according to how well it performed on the training set. Figure 18.33 shows how the algorithm works conceptually. There are many variants of the basic boosting idea, with different ways of adjusting the weights and combining the hypotheses. One specific algorithm, called ADABOOST, is shown in Figure 18.34. ADABOOST has a very important property: if the input learning algorithm  $L$  is a **weak learning** algorithm—which

WEAK LEARNING

<sup>14</sup> For learning algorithms in which this is not possible, one can instead create a **replicated training set** where the  $j$ th example appears  $w_j$  times, using randomization to handle fractional weights.



means that  $L$  always returns a hypothesis with accuracy on the training set that is slightly better than random guessing (i.e.,  $50\% + \epsilon$  for Boolean classification)—then ADABOOST will return a hypothesis that *classifies the training data perfectly* for large enough  $K$ . Thus, the algorithm *boosts* the accuracy of the original learning algorithm on the training data. This result holds no matter how inexpressive the original hypothesis space and no matter how complex the function being learned.

#### DECISION STUMP

Let us see how well boosting does on the restaurant data. We will choose as our original hypothesis space the class of **decision stumps**, which are decision trees with just one test, at the root. The lower curve in Figure 18.35(a) shows that unboosted decision stumps are not very effective for this data set, reaching a prediction performance of only 81% on 100 training examples. When boosting is applied (with  $K = 5$ ), the performance is better, reaching 93% after 100 examples.

An interesting thing happens as the ensemble size  $K$  increases. Figure 18.35(b) shows the training set performance (on 100 examples) as a function of  $K$ . Notice that the error reaches zero when  $K$  is 20; that is, a weighted-majority combination of 20 decision stumps suffices to fit the 100 examples exactly. As more stumps are added to the ensemble, the error remains at zero. The graph also shows that *the test set performance continues to increase long after the training set error has reached zero*. At  $K = 20$ , the test performance is 0.95 (or 0.05 error), and the performance increases to 0.98 as late as  $K = 137$ , before gradually dropping to 0.95.

This finding, which is quite robust across data sets and hypothesis spaces, came as quite a surprise when it was first noticed. Ockham's razor tells us not to make hypotheses more

