

# 3

## *The Bayesian Network Representation*



Our goal is to represent a joint distribution  $P$  over some set of random variables  $\mathcal{X} = \{X_1, \dots, X_n\}$ . Even in the simplest case where these variables are binary-valued, a joint distribution requires the specification of  $2^n - 1$  numbers — the probabilities of the  $2^n$  different assignments of values  $x_1, \dots, x_n$ . For all but the smallest  $n$ , the explicit representation of the joint distribution is unmanageable from every perspective. Computationally, it is very expensive to manipulate and generally too large to store in memory. Cognitively, it is impossible to acquire so many numbers from a human experts; moreover, the numbers are very small and do not correspond to events that people can reasonably contemplate. Statistically, if we want to learn the distribution from data, we would need ridiculously large amounts of data to estimate this many parameters robustly. These problems were the main barrier to the adoption of probabilistic methods for expert systems until the development of the methodologies described in this book.

In this chapter, we first show how independence properties in the distribution can be used to represent such high-dimensional distributions much more compactly. We then show how a combinatorial data structure — a directed acyclic graph — can provide us with a general-purpose modeling language for exploiting this type of structure in our representation.

### 3.1 Exploiting Independence Properties

The compact representations we explore in this chapter are based on two key ideas: the representation of independence properties of the distribution, and the use of an alternative parameterization that allows us to exploit these finer-grained independencies.

#### 3.1.1 Independent Random Variables

To motivate our discussion, consider a simple setting where we know that each  $X_i$  represents the outcome of a toss of coin  $i$ . In this case, we typically assume that the different coin tosses are marginally independent (definition 2.4), so that our distribution  $P$  will satisfy  $(X_i \perp X_j)$  for any  $i, j$ . More generally (strictly more generally — see exercise 3.1), we assume that the distribution satisfies  $(X \perp Y)$  for any disjoint subsets of the variables  $X$  and  $Y$ . Therefore, we have that:

$$P(X_1, \dots, X_n) = P(X_1)P(X_2) \cdots P(X_n).$$

parameters

If we use the standard parameterization of the joint distribution, this independence structure is obscured, and the representation of the distribution requires  $2^n$  parameters. However, we can use a more natural set of parameters for specifying this distribution: If  $\theta_i$  is the probability with which coin  $i$  lands heads, the joint distribution  $P$  can be specified using the  $n$  parameters  $\theta_1, \dots, \theta_n$ . These parameters implicitly specify the  $2^n$  probabilities in the joint distribution. For example, the probability that all of the coin tosses land heads is simply  $\theta_1 \cdot \theta_2 \cdot \dots \cdot \theta_n$ . More generally, letting  $\theta_{x_i} = \theta_i$  when  $x_i = x_i^1$  and  $\theta_{x_i} = 1 - \theta_i$  when  $x_i = x_i^0$ , we can define:

$$P(x_1, \dots, x_n) = \prod_i \theta_{x_i}. \quad (3.1)$$

independent parameters

This representation is limited, and there are many distributions that we cannot capture by choosing values for  $\theta_1, \dots, \theta_n$ . This fact is obvious not only from intuition, but also from a somewhat more formal perspective. The space of all joint distributions is a  $2^n - 1$  dimensional subspace of  $\mathbb{R}^{2^n}$  — the set  $\{(p_1, \dots, p_{2^n}) \in \mathbb{R}^{2^n} : p_1 + \dots + p_{2^n} = 1\}$ . On the other hand, the space of all joint distributions specified in a factorized way as in equation (3.1) is an  $n$ -dimensional manifold in  $\mathbb{R}^{2^n}$ .

A key concept here is the notion of *independent parameters* — parameters whose values are not determined by others. For example, when specifying an arbitrary multinomial distribution over a  $k$  dimensional space, we have  $k - 1$  independent parameters: the last probability is fully determined by the first  $k - 1$ . In the case where we have an arbitrary joint distribution over  $n$  binary random variables, the number of independent parameters is  $2^n - 1$ . On the other hand, the number of independent parameters for distributions represented as  $n$  independent binomial coin tosses is  $n$ . Therefore, the two spaces of distributions cannot be the same. (While this argument might seem trivial in this simple case, it turns out to be an important tool for comparing the expressive power of different representations.)

As this simple example shows, certain families of distributions — in this case, the distributions generated by  $n$  independent random variables — permit an alternative parameterization that is substantially more compact than the naive representation as an explicit joint distribution. Of course, in most real-world applications, the random variables are not marginally independent. However, a generalization of this approach will be the basis for our solution.

### 3.1.2 The Conditional Parameterization

Let us begin with a simple example that illustrates the basic intuition. Consider the problem faced by a company trying to hire a recent college graduate. The company's goal is to hire intelligent employees, but there is no way to test intelligence directly. However, the company has access to the student's SAT scores, which are informative but not fully indicative. Thus, our probability space is induced by the two random variables *Intelligence* ( $I$ ) and *SAT* ( $S$ ). For simplicity, we assume that each of these takes two values:  $Val(I) = \{i^1, i^0\}$ , which represent the values high intelligence ( $i^1$ ) and low intelligence ( $i^0$ ); similarly  $Val(S) = \{s^1, s^0\}$ , which also represent the values *high* (score) and *low* (score), respectively.

Thus, our joint distribution in this case has four entries. For example, one possible joint

distribution  $P$  would be

$I$	$S$	$P(I, S)$
$i^0$	$s^0$	0.665
$i^0$	$s^1$	0.035
$i^1$	$s^0$	0.06
$i^1$	$s^1$	0.24

(3.2)

There is, however, an alternative, and even more natural way of representing the same joint distribution. Using the chain rule of conditional probabilities (see equation (2.5)), we have that

$$P(I, S) = P(I)P(S | I).$$

Intuitively, we are representing the process in a way that is more compatible with causality. Various factors (genetics, upbringing, ...) first determined (stochastically) the student's intelligence. His performance on the SAT is determined (stochastically) by his intelligence. We note that the models we construct are not required to follow causal intuitions, but they often do. We return to this issue later on.

From a mathematical perspective, this equation leads to the following alternative way of representing the joint distribution. Instead of specifying the various joint entries  $P(I, S)$ , we would specify it in the form of  $P(I)$  and  $P(S | I)$ . Thus, for example, we can represent the joint distribution of equation (3.2) using the following two tables, one representing the *prior distribution* over  $I$  and the other the *conditional probability distribution (CPD)* of  $S$  given  $I$ :

prior distribution  
CPD

$i^0$	$i^1$	$I$	$s^0$	$s^1$
0.7	0.3	$i^0$	0.95	0.05
		$i^1$	0.2	0.8

(3.3)

The CPD  $P(S | I)$  represents the probability that the student will succeed on his SATs in the two possible cases: the case where the student's intelligence is low, and the case where it is high. The CPD asserts that a student of low intelligence is extremely unlikely to get a high SAT score ( $P(s^1 | i^0) = 0.05$ ); on the other hand, a student of high intelligence is likely, but far from certain, to get a high SAT score ( $P(s^1 | i^1) = 0.8$ ).

It is instructive to consider how we could parameterize this alternative representation. Here, we are using three binomial distributions, one for  $P(I)$ , and two for  $P(S | i^0)$  and  $P(S | i^1)$ . Hence, we can parameterize this representation using three independent parameters, say  $\theta_{i^1}$ ,  $\theta_{s^1|i^1}$ , and  $\theta_{s^1|i^0}$ . Our representation of the joint distribution as a four-outcome multinomial also required three parameters. Thus, although the conditional representation is more natural than the explicit representation of the joint, it is not more compact. However, as we will soon see, the conditional parameterization provides a basis for our compact representations of more complex distributions.

Although we will only define Bayesian networks formally in section 3.2.2, it is instructive to see how this example would be represented as one. The Bayesian network, as shown in figure 3.1a, would have a node for each of the two random variables  $I$  and  $S$ , with an edge from  $I$  to  $S$  representing the direction of the dependence in this model.

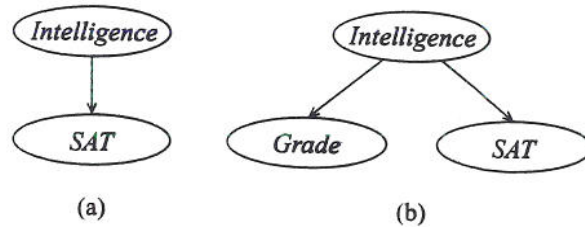


Figure 3.1 Simple Bayesian networks for the student example

### 3.1.3 The Naive Bayes Model

We now describe perhaps the simplest example where a conditional parameterization is combined with conditional independence assumptions to produce a very compact representation of a high-dimensional probability distribution. Importantly, unlike the previous example of fully independent random variables, none of the variables in this distribution are (marginally) independent.

#### 3.1.3.1 The Student Example

Elaborating our example, we now assume that the company also has access to the student's grade  $G$  in some course. In this case, our probability space is the joint distribution over the three relevant random variables  $I$ ,  $S$ , and  $G$ . Assuming that  $I$  and  $S$  are as before, and that  $G$  takes on three values  $g^1, g^2, g^3$ , representing the grades  $A, B$ , and  $C$ , respectively, then the joint distribution has twelve entries.

Before we even consider the specific numerical aspects of our distribution  $P$  in this example, we can see that independence does not help us: for any reasonable  $P$ , there are no independencies that hold. The student's intelligence is clearly correlated both with his SAT score and with his grade. The SAT score and grade are also not independent: if we condition on the fact that the student received a high score on his SAT, the chances that he gets a high grade in his class are also likely to increase. Thus, we may assume that, for our particular distribution  $P$   $P(g^1 | s^1) > P(g^1 | s^0)$ .

However, it is quite plausible that our distribution  $P$  in this case satisfies a conditional independence property. If we know that the student has high intelligence, a high grade on the SAT no longer gives us information about the student's performance in the class. More formally:

$$P(g | i^1, s^1) = P(g | i^1).$$

More generally, we may well assume that

$$P \models (S \perp G | I). \tag{3.4}$$

Note that this independence statement holds only if we assume that the student's intelligence is the only reason why his grade and SAT score might be correlated. In other words, it assumes that there are no correlations due to other factors, such as the student's ability to take timed exams. These assumptions are also not "true" in any formal sense of the word, and they are often only approximations of our true beliefs. (See box 3.C for some further discussion.)

As in the case of marginal independence, conditional independence allows us to provide a compact specification of the joint distribution. Again, the compact representation is based on a very natural alternative parameterization. By simple probabilistic reasoning (as in equation (2.5)), we have that

$$P(I, S, G) = P(S, G | I)P(I).$$

But now, the conditional independence assumption of equation (3.4) implies that

$$P(S, G | I) = P(S | I)P(G | I).$$

Hence, we have that

$$P(I, S, G) = P(S | I)P(G | I)P(I). \quad (3.5)$$

Thus, we have factorized the joint distribution  $P(I, S, G)$  as a product of three conditional probability distributions (CPDs). This factorization immediately leads us to the desired alternative parameterization. In order to specify fully a joint distribution satisfying equation (3.4), we need the following three CPDs:  $P(I)$ ,  $P(S | I)$ , and  $P(G | I)$ . The first two might be the same as in equation (3.3). The latter might be

$I$	$g^1$	$g^2$	$g^3$
$i^0$	0.2	0.34	0.46
$i^1$	0.74	0.17	0.09

Together, these three CPDs fully specify the joint distribution (assuming the conditional independence of equation (3.4)). For example,

$$\begin{aligned} P(i^1, s^1, g^2) &= P(i^1)P(s^1 | i^1)P(g^2 | i^1) \\ &= 0.3 \cdot 0.8 \cdot 0.17 = 0.0408. \end{aligned}$$

Once again, we note that this probabilistic model would be represented using the Bayesian network shown in figure 3.1b.

In this case, the alternative parameterization is more compact than the joint. We now have three binomial distributions —  $P(I)$ ,  $P(S | i^1)$  and  $P(S | i^0)$ , and two three-valued multinomial distributions —  $P(G | i^1)$  and  $P(G | i^0)$ . Each of the binomials requires one independent parameter, and each three-valued multinomial requires two independent parameters, for a total of seven. By contrast, our joint distribution has twelve entries, so that eleven independent parameters are required to specify an arbitrary joint distribution over these three variables.

It is important to note another advantage of this way of representing the joint: modularity. When we added the new variable  $G$ , the joint distribution changed entirely. Had we used the explicit representation of the joint, we would have had to write down twelve new numbers. In the factored representation, we could reuse our local probability models for the variables  $I$  and  $S$ , and specify only the probability model for  $G$  — the CPD  $P(G | I)$ . This property will turn out to be invaluable in modeling real-world systems.

### 3.1.3.2 The General Model

naive Bayes

This example is an instance of a much more general model commonly called the *naive Bayes*

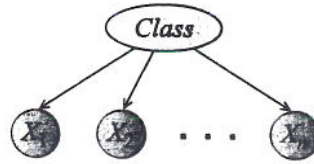


Figure 3.2 The Bayesian network graph for a naive Bayes model

model (also known as the *Idiot Bayes model*). The naive Bayes model assumes that instances fall into one of a number of mutually exclusive and exhaustive *classes*. Thus, we have a class variable  $C$  that takes on values in some set  $\{c^1, \dots, c^k\}$ . In our example, the class variable is the student's intelligence  $I$ , and there are two classes of instances — students with high intelligence and students with low intelligence.

features

The model also includes some number of *features*  $X_1, \dots, X_k$  whose values are typically observed. The *naive Bayes assumption* is that the features are conditionally independent given the instance's class. In other words, within each class of instances, the different properties can be determined independently. More formally, we have that

$$(X_i \perp X_{-i} \mid C) \quad \text{for all } i, \quad (3.6)$$

where  $X_{-i} = \{X_1, \dots, X_k\} - \{X_i\}$ . This model can be represented using the Bayesian network of figure 3.2. In this example, and later on in the book, we use a darker oval to represent variables that are always observed when the network is used.

factorization

Based on these independence assumptions, we can show that the model *factorizes* as:

$$P(C, X_1, \dots, X_n) = P(C) \prod_{i=1}^n P(X_i \mid C). \quad (3.7)$$

(See exercise 3.2.) Thus, in this model, we can represent the joint distribution using a small set of factors: a prior distribution  $P(C)$ , specifying how likely an instance is to belong to different classes a priori, and a set of CPDs  $P(X_j \mid C)$ , one for each of the  $n$  finding variables. These factors can be encoded using a very small number of parameters. For example, if all of the variables are binary, the number of independent parameters required to specify the distribution is  $2n + 1$  (see exercise 3.6). Thus, the number of parameters is linear in the number of variables, as opposed to exponential for the explicit representation of the joint.

classification

**Box 3.A — Concept: The Naive Bayes Model.** *The naive Bayes model, despite the strong assumptions that it makes, is often used in practice, because of its simplicity and the small number of parameters required. The model is generally used for classification — deciding, based on the values of the evidence variables for a given instance, the class to which the instance is most likely to belong. We might also want to compute our confidence in this decision, that is, the extent to which our model favors one class  $c^1$  over another  $c^2$ . Both queries can be addressed by the following ratio:*

$$\frac{P(C = c^1 \mid x_1, \dots, x_n)}{P(C = c^2 \mid x_1, \dots, x_n)} = \frac{P(C = c^1)}{P(C = c^2)} \prod_{i=1}^n \frac{P(x_i \mid C = c^1)}{P(x_i \mid C = c^2)}; \quad (3.8)$$

see exercise 3.2). This formula is very natural, since it computes the posterior probability ratio of  $c^1$  versus  $c^2$  as a product of their prior probability ratio (the first term), multiplied by a set of terms  $\frac{P(x_i|C=c^1)}{P(x_i|C=c^2)}$  that measure the relative support of the finding  $x_i$  for the two classes.

medical diagnosis

This model was used in the early days of medical diagnosis, where the different values of the class variable represented different diseases that the patient could have. The evidence variables represented different symptoms, test results, and the like. Note that the model makes several strong assumptions that are not generally true, specifically that the patient can have at most one disease, and that, given the patient's disease, the presence or absence of different symptoms, and the values of different tests, are all independent. This model was used for medical diagnosis because the small number of interpretable parameters made it easy to elicit from experts. For example, it is quite natural to ask of an expert physician what the probability is that a patient with pneumonia has high fever. Indeed, several early medical diagnosis systems were based on this technology, and some were shown to provide better diagnoses than those made by expert physicians.

However, later experience showed that the strong assumptions underlying this model decrease its diagnostic accuracy. In particular, the model tends to overestimate the impact of certain evidence by "overcounting" it. For example, both hypertension (high blood pressure) and obesity are strong indicators of heart disease. However, because these two symptoms are themselves highly correlated, equation (3.8), which contains a multiplicative term for each of them, double-counts the evidence they provide about the disease. Indeed, some studies showed that the diagnostic performance of a naive Bayes model degraded as the number of features increased; this degradation was often traced to violations of the strong conditional independence assumption. This phenomenon led to the use of more complex Bayesian networks, with more realistic independence assumptions, for this application (see box 3.D).

Nevertheless, the naive Bayes model is still useful in a variety of applications, particularly in the context of models learned from data in domains with a large number of features and a relatively small number of instances, such as classifying documents into topics using the words in the documents as features; see box 17.E).

## 3.2 Bayesian Networks

Bayesian networks build on the same intuitions as the naive Bayes model by exploiting conditional independence properties of the distribution in order to allow a compact and natural representation. However, they are not restricted to representing distributions satisfying the strong independence assumptions implicit in the naive Bayes model. They allow us the flexibility to tailor our representation of the distribution to the independence properties that appear reasonable in the current setting.

The core of the Bayesian network representation is a directed acyclic graph (DAG)  $\mathcal{G}$ , whose nodes are the random variables in our domain and whose edges correspond, intuitively, to direct influence of one node on another. This graph  $\mathcal{G}$  can be viewed in two very different ways:



- as a data structure that provides the skeleton for representing a joint distribution compactly in a factorized way;

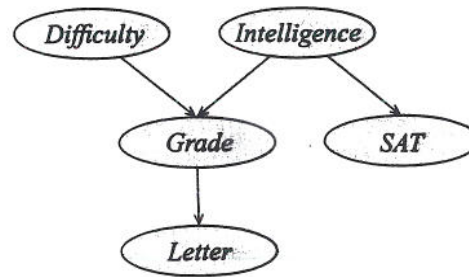


Figure 3.3 The Bayesian Network graph for the Student example

- as a compact representation for a set of conditional independence assumptions about a distribution.

As we will see, these two views are, in a strong sense, equivalent.

### 3.2.1 The Student Example Revisited

We begin our discussion with a simple toy example, which will accompany us, in various versions, throughout much of this book.

#### 3.2.1.1 The Model

Consider our student from before, but now consider a slightly more complex scenario. The student's grade, in this case, depends not only on his intelligence but also on the difficulty of the course, represented by a random variable  $D$  whose domain is  $Val(D) = \{easy, hard\}$ . Our student asks his professor for a recommendation letter. The professor is absentminded and never remembers the names of her students. She can only look at his grade, and she writes her letter for him based on that information alone. The quality of her letter is a random variable  $L$ , whose domain is  $Val(L) = \{strong, weak\}$ . The actual quality of the letter depends stochastically on the grade. (It can vary depending on how stressed the professor is and the quality of the coffee she had that morning.)

We therefore have five random variables in this domain: the student's intelligence ( $I$ ), the course difficulty ( $D$ ), the grade ( $G$ ), the student's SAT score ( $S$ ), and the quality of the recommendation letter ( $L$ ). All of the variables except  $G$  are binary-valued, and  $G$  is ternary-valued. Hence, the joint distribution has 48 entries.

As we saw in our simple illustrations of figure 3.1, a Bayesian network is represented using a directed graph whose nodes represent the random variables and whose edges represent direct influence of one variable on another. We can view the graph as encoding a generative sampling process executed by nature, where the value for each variable is selected by nature using a distribution that depends only on its parents. In other words, each variable is a stochastic function of its parents.

Based on this intuition, perhaps the most natural network structure for the distribution in this example is the one presented in figure 3.3. The edges encode our intuition about



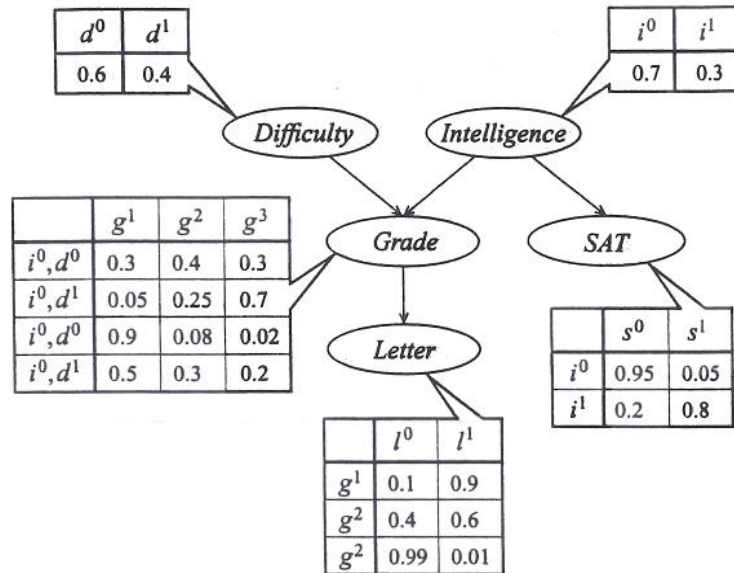


Figure 3.4 Student Bayesian network  $B^{student}$  with CPDs

the way the world works. The course difficulty and the student's intelligence are determined independently, and before any of the variables in the model. The student's grade depends on both of these factors. The student's SAT score depends only on his intelligence. The quality of the professor's recommendation letter depends (by assumption) only on the student's grade in the class. Intuitively, each variable in the model depends directly only on its parents in the network. We formalize this intuition later.

The second component of the Bayesian network representation is a set of *local probability models* that represent the nature of the dependence of each variable on its parents. One such model,  $P(I)$ , represents the distribution in the population of intelligent versus less intelligent student. Another,  $P(D)$ , represents the distribution of difficult and easy classes. The distribution over the student's grade is a conditional distribution  $P(G | I, D)$ . It specifies the distribution over the student's grade, inasmuch as it depends on the student's intelligence and the course difficulty. Specifically, we would have a different distribution for each assignment of values  $i, d$ . For example, we might believe that a smart student in an easy class is 90 percent likely to get an A, 8 percent likely to get a B, and 2 percent likely to get a C. Conversely, a smart student in a hard class may only be 30 percent likely to get an A. In general, each variable  $X$  in the model is associated with a *conditional probability distribution (CPD)* that specifies a distribution over the values of  $X$  given each possible joint assignment of values to its parents in the model. For a node with no parents, the CPD is conditioned on the empty set of variables. Hence, the CPD turns into a marginal distribution, such as  $P(D)$  or  $P(I)$ . One possible choice of CPDs for this domain is shown in figure 3.4. The network structure together with its CPDs is a *Bayesian network*  $B$ ; we use  $B^{student}$  to refer to the Bayesian network for our student example.

local probability  
model

CPD

How do we use this data structure to specify the joint distribution? Consider some particular state in this space, for example,  $i^1, d^0, g^2, s^1, l^0$ . Intuitively, the probability of this event can be computed from the probabilities of the basic events that comprise it: the probability that the student is intelligent; the probability that the course is easy; the probability that a smart student gets a B in an easy class; the probability that a smart student gets a high score on his SAT; and the probability that a student who got a B in the class gets a weak letter. The total probability of this state is:

$$\begin{aligned} P(i^1, d^0, g^2, s^1, l^0) &= P(i^1)P(d^0)P(g^2 | i^1, d^0)P(s^1 | i^1)P(l^0 | g^2) \\ &= 0.3 \cdot 0.6 \cdot 0.08 \cdot 0.8 \cdot 0.4 = 0.004608. \end{aligned}$$

Clearly, we can use the same process for any state in the joint probability space. In general, we will have that

$$P(I, D, G, S, L) = P(I)P(D)P(G | I, D)P(S | I)P(L | G). \quad (3.9)$$

chain rule for  
Bayesian  
networks

This equation is our first example of the *chain rule for Bayesian networks* which we will define in a general setting in section 3.2.3.2.

### 3.2.1.2 Reasoning Patterns

A joint distribution  $P_{\mathcal{B}}$  specifies (albeit implicitly) the probability  $P_{\mathcal{B}}(Y = y | E = e)$  of any event  $y$  given any observations  $e$ , as discussed in section 2.1.3.3: We condition the joint distribution on the event  $E = e$  by eliminating the entries in the joint inconsistent with our observation  $e$ , and renormalizing the resulting entries to sum to 1; we compute the probability of the event  $y$  by summing the probabilities of all of the entries in the resulting posterior distribution that are consistent with  $y$ . To illustrate this process, let us consider our  $\mathcal{B}^{student}$  network and see how the probabilities of various events change as evidence is obtained.

Consider a particular student, George, about whom we would like to reason using our model. We might ask how likely George is to get a strong recommendation ( $l^1$ ) from his professor in Econ101. Knowing nothing else about George or Econ101, this probability is about 50.2 percent. More precisely, let  $P_{\mathcal{B}^{student}}$  be the joint distribution defined by the preceding BN; then we have that  $P_{\mathcal{B}^{student}}(l^1) \approx 0.502$ . We now find out that George is not so intelligent ( $i^0$ ); the probability that he gets a strong letter from the professor of Econ101 goes down to around 38.9 percent; that is,  $P_{\mathcal{B}^{student}}(l^1 | i^0) \approx 0.389$ . We now further discover that Econ101 is an easy class ( $d^0$ ). The probability that George gets a strong letter from the professor is now  $P_{\mathcal{B}^{student}}(l^1 | i^0, d^0) \approx 0.51$ . Queries such as these, where we predict the “downstream” effects of various factors (such as George’s intelligence), are instances of *causal reasoning* or *prediction*.

causal reasoning

Now, consider a recruiter for Acme Consulting, trying to decide whether to hire George based on our previous model. A priori, the recruiter believes that George is 30 percent likely to be intelligent. He obtains George’s grade record for a particular class Econ101 and sees that George received a C in the class ( $g^3$ ). His probability that George has high intelligence goes down significantly, to about 7.9 percent; that is,  $P_{\mathcal{B}^{student}}(i^1 | g^3) \approx 0.079$ . We note that the probability that the class is a difficult one also goes up, from 40 percent to 62.9 percent.

Now, assume that the recruiter fortunately (for George) lost George’s transcript, and has only the recommendation letter from George’s professor in Econ101, which (not surprisingly)

evidential  
reasoning

weak. The probability that George has high intelligence still goes down, but only to 14 percent:  $P_{\text{student}}(i^1 | l^0) \approx 0.14$ . Note that if the recruiter has both the grade and the letter, we have the same probability as if he had only the grade:  $P_{\text{student}}(i^1 | g^3, l^0) \approx 0.079$ ; we will revisit this issue. Queries such as this, where we reason from effects to causes, are instances of *evidential reasoning* or *explanation*.

Finally, George submits his SAT scores to the recruiter, and astonishingly, his SAT score is high. The probability that George has high intelligence goes up dramatically, from 7.9 percent to 57.8 percent:  $P_{\text{student}}(i^1 | g^3, s^1) \approx 0.578$ . Intuitively, the reason that the high SAT score outweighs the poor grade is that students with low intelligence are extremely unlikely to get good scores on their SAT, whereas students with high intelligence can still get C's. However, smart students are much more likely to get C's in hard classes. Indeed, we see that the probability that Econ101 is a difficult class goes up from the 62.9 percent we saw before to around 76 percent.

This last pattern of reasoning is a particularly interesting one. The information about the SAT gave us information about the student's intelligence, which, in conjunction with the student's grade in the course, told us something about the difficulty of the course. In effect, we have one causal factor for the *Grade* variable — *Intelligence* — giving us information about another — *Difficulty*.

explaining away

intercausal  
reasoning

Let us examine this pattern in its pure form. As we said,  $P_{\text{student}}(i^1 | g^3) \approx 0.079$ . On the other hand, if we now discover that Econ101 is a hard class, we have that  $P_{\text{student}}(i^1 | g^3, d^1) \approx 0.11$ . In effect, we have provided at least a partial explanation for George's grade in Econ101. To take an even more striking example, if George gets a B in Econ 101, we have that  $P_{\text{student}}(i^1 | g^2) \approx 0.175$ . On the other hand, if Econ101 is a hard class, we get  $P_{\text{student}}(i^1 | g^2, d^1) \approx 0.34$ . In effect we have *explained away* the poor grade via the difficulty of the class. **Explaining away is an instance of a general reasoning pattern called *intercausal reasoning*, where different causes of the same effect can interact. This type of reasoning is a very common pattern in human reasoning.** For example, when we have fever and a sore throat, and are concerned about mononucleosis, we are greatly relieved to be told we have the flu. Clearly, having the flu does not prohibit us from having mononucleosis. Yet, having the flu provides an alternative explanation of our symptoms, thereby reducing substantially the probability of mononucleosis.

This intuition of providing an alternative explanation for the evidence can be made very precise. As shown in exercise 3.3, if the flu deterministically causes the symptoms, the probability of mononucleosis goes down to its prior probability (the one prior to the observations of any symptoms). On the other hand, if the flu might occur without causing these symptoms, the probability of mononucleosis goes down, but it still remains somewhat higher than its base level. Explaining away, however, is not the only form of intercausal reasoning. The influence can go in any direction. Consider, for example, a situation where someone is found dead and may have been murdered. The probabilities that a suspect has motive and opportunity both go up. If we now discover that the suspect has motive, the probability that he has opportunity goes up. (See exercise 3.4.)

It is important to emphasize that, although our explanations used intuitive concepts such as cause and evidence, there is nothing mysterious about the probability computations we performed. They can be replicated simply by generating the joint distribution, as defined in equation (3.9), and computing the probabilities of the various events directly from that.

### 3.2.2 Basic Independencies in Bayesian Networks

As we discussed, a Bayesian network graph  $\mathcal{G}$  can be viewed in two ways. In the previous section, we showed, by example, how it can be used as a skeleton data structure to which we can attach local probability models that together define a joint distribution. In this section, we provide a formal semantics for a Bayesian network, starting from the perspective that the graph encodes a set of conditional independence assumptions. We begin by understanding, intuitively, the basic conditional independence assumptions that we want a directed graph to encode. We then formalize these desired assumptions in a definition.

#### 3.2.2.1 Independencies in the Student Example

In the Student example, we used the intuition that edges represent direct dependence. For example, we made intuitive statements such as “the professor’s recommendation letter depends only on the student’s grade in the class”; this statement was encoded in the graph by the fact that there are no direct edges into the  $L$  node except from  $G$ . This intuition, that “a node depends directly only on its parents,” lies at the heart of the semantics of Bayesian networks.

We give formal semantics to this assertion using conditional independence statements. For example, the previous assertion can be stated formally as the assumption that  $L$  is conditionally independent of all other nodes in the network given its parent  $G$ :

$$(L \perp I, D, S \mid G). \quad (3.10)$$

In other words, once we know the student’s grade, our beliefs about the quality of his recommendation letter are not influenced by information about any other variable. Similarly, to formalize our intuition that the student’s SAT score depends only on his intelligence, we can say that  $S$  is conditionally independent of all other nodes in the network given its parent  $I$ :

$$(S \perp D, G, L \mid I). \quad (3.11)$$

Now, let us consider the  $G$  node. Following the pattern blindly, we may be tempted to assert that  $G$  is conditionally independent of all other variables in the network given its parents. However, this assumption is false both at an intuitive level and for the specific example distribution we used earlier. Assume, for example, that we condition on  $i^1, d^1$ ; that is, we have a smart student in a difficult class. In this setting, is  $G$  independent of  $L$ ? Clearly, the answer is no: if we observe  $l^1$  (the student got a strong letter), then our probability in  $g^1$  (the student received an A in the course) should go up; that is, we would expect

$$P(g^1 \mid i^1, d^1, l^1) > P(g^1 \mid i^1, d^1).$$

Indeed, if we examine our distribution, the latter probability is 0.5 (as specified in the CPD), whereas the former is a much higher 0.712.

Thus, we see that we do not expect a node to be conditionally independent of all other nodes given its parents. In particular, even given its parents, it can still depend on its descendants. Can it depend on other nodes? For example, do we expect  $G$  to depend on  $S$  given  $I$  and  $D$ ? Intuitively, the answer is no. Once we know, say, that the student has high intelligence, his SAT score gives us no additional information that is relevant toward predicting his grade. Thus, we

would want the property that:

$$(G \perp S \mid I, D). \quad (3.12)$$

It remains only to consider the variables  $I$  and  $D$ , which have no parents in the graph. Thus, in our search for independencies given a node's parents, we are now looking for marginal independencies. As the preceding discussion shows, in our distribution  $P_{\mathcal{B}_{student}}$ ,  $I$  is not independent of its descendants  $G$ ,  $L$ , or  $S$ . Indeed, the only nondescendant of  $I$  is  $D$ . Indeed, we assumed implicitly that *Intelligence* and *Difficulty* are independent. Thus, we expect that:

$$(I \perp D). \quad (3.13)$$

This analysis might seem somewhat surprising in light of our earlier examples, where learning something about the course difficulty drastically changed our beliefs about the student's intelligence. In that situation, however, we were reasoning in the presence of information about the student's grade. In other words, we were demonstrating the dependence of  $I$  and  $D$  given  $G$ . This phenomenon is a very important one, and we will return to it.

For the variable  $D$ , both  $I$  and  $S$  are nondescendants. Recall that, if  $(I \perp D)$  then  $(D \perp I)$ . The variable  $S$  increases our beliefs in the student's intelligence, but knowing that the student is smart (or not) does not influence our beliefs in the difficulty of the course. Thus, we have that

$$(D \perp I, S). \quad (3.14)$$

We can see a pattern emerging. Our intuition tells us that the parents of a variable "shield" it from probabilistic influence that is causal in nature. In other words, once I know the value of the parents, no information relating directly or indirectly to its parents or other ancestors can influence my beliefs about it. However, information about its descendants *can* change my beliefs about it, via an evidential reasoning process.

### 3.2.2.2 Bayesian Network Semantics

We are now ready to provide the formal definition of the semantics of a Bayesian network structure. We would like the formal definition to match the intuitions developed in our example.

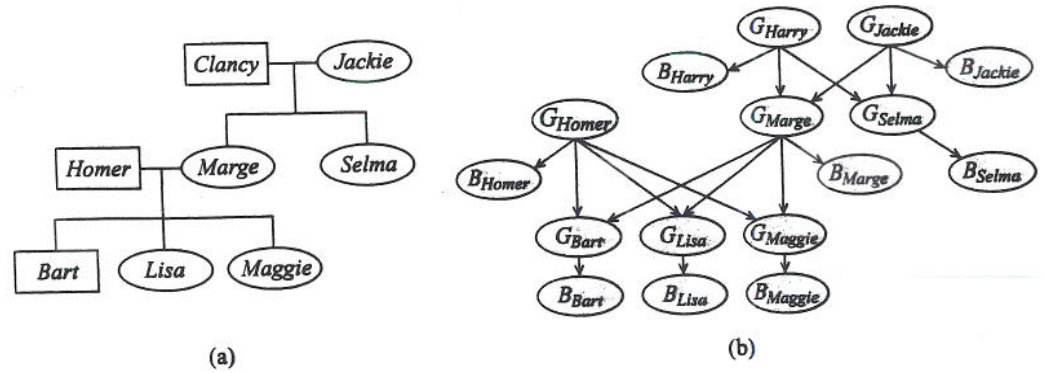
**Definition 3.1**  
Bayesian network  
structure  
local  
independencies

A Bayesian network structure  $\mathcal{G}$  is a directed acyclic graph whose nodes represent random variables  $X_1, \dots, X_n$ . Let  $\text{Pa}_{X_i}^{\mathcal{G}}$  denote the parents of  $X_i$  in  $\mathcal{G}$ , and  $\text{NonDescendants}_{X_i}$  denote the variables in the graph that are not descendants of  $X_i$ . Then  $\mathcal{G}$  encodes the following set of conditional independence assumptions, called the local independencies, and denoted by  $\mathcal{I}_{\ell}(\mathcal{G})$ :

$$\text{For each variable } X_i: (X_i \perp \text{NonDescendants}_{X_i} \mid \text{Pa}_{X_i}^{\mathcal{G}}). \quad \blacksquare$$

In other words, the local independencies state that each node  $X_i$  is conditionally independent of its nondescendants given its parents.

Returning to the Student network  $G_{student}$ , the local Markov independencies are precisely the ones dictated by our intuition, and specified in equation (3.10) – equation (3.14).



**Figure 3.B.1 — Modeling Genetic Inheritance** (a) A small family tree. (b) A simple BN for genetic inheritance in this domain. The  $G$  variables represent a person's genotype, and the  $B$  variables the result of a blood-type test.

**Box 3.B — Case Study: The Genetics Example.** One of the very earliest uses of a Bayesian network model (long before the general framework was defined) is in the area of genetic pedigrees. In this setting, the local independencies are particularly intuitive. In this application, we want to model the transmission of a certain property, say blood type, from parent to child. The blood type of a person is an observable quantity that depends on her genetic makeup. Such properties are called phenotypes. The genetic makeup of a person is called genotype.

To model this scenario properly, we need to introduce some background on genetics. The human genetic material consists of 22 pairs of autosomal chromosomes and a pair of the sex chromosomes ( $X$  and  $Y$ ). Each chromosome contains a set of genetic material, consisting (among other things) of genes that determine a person's properties. A region of the chromosome that is of interest is called a locus; a locus can have several variants, called alleles.

For concreteness, we focus on autosomal chromosome pairs. In each autosomal pair, one chromosome is the paternal chromosome, inherited from the father, and the other is the maternal chromosome, inherited from the mother. For genes in an autosomal pair, a person has two copies of the gene, one on each copy of the chromosome. Thus, one of the gene's alleles is inherited from the person's mother, and the other from the person's father. For example, the region containing the gene that encodes a person's blood type is a locus. This gene comes in three variants, or alleles:  $A$ ,  $B$ , and  $O$ . Thus, a person's genotype is denoted by an ordered pair, such as  $\langle A, B \rangle$ ; with three choices for each entry in the pair, there are 9 possible genotypes. The blood type phenotype is a function of both copies of the gene. For example, if the person has an  $A$  allele and an  $O$  allele, her observed blood type is "A." If she has two  $O$  alleles, her observed blood type is "O."

To represent this domain, we would have, for each person, two variables: one representing the person's genotype, and the other her phenotype. We use the name  $G(p)$  to represent person  $p$ 's genotype, and  $B(p)$  to represent her blood type.

In this example, the independence assumptions arise immediately from the biology. Since the

blood type is a function of the genotype, once we know the genotype of a person, additional evidence about other members of the family will not provide new information about the blood type. Similarly, the process of genetic inheritance implies independence assumption. Once we know the genotype of both parents, we know what each of them can pass on to the offspring. Thus, learning new information about ancestors (or nondescendants) does not provide new information about the genotype of the offspring. These are precisely the local independencies in the resulting network structure, shown for a simple family tree in figure 3.B.1. The intuition here is clear; for example, Bart's blood type is correlated with that of his aunt Selma, but once we know Homer's and Marge's genotype, the two become independent.

To define the probabilistic model fully, we need to specify the CPDs. There are three types of CPDs in this model:

- The penetrance model  $P(B(c) | G(c))$ , which describes the probability of different variants of a particular phenotype (say different blood types) given the person's genotype. In the case of the blood type, this CPD is a deterministic function, but in other cases, the dependence can be more complex.
- The transmission model  $P(G(c) | G(p), G(m))$ , where  $c$  is a person and  $p, m$  her father and mother, respectively. Each parent is equally likely to transmit either of his or her two alleles to the child.
- Genotype priors  $P(G(c))$ , used when person  $c$  has no parents in the pedigree. These are the general genotype frequencies within the population.

Our discussion of blood type is simplified for several reasons. First, some phenotypes, such as late-onset diseases, are not a deterministic function of the genotype. Rather, an individual with a particular genotype might be more likely to have the disease than an individual with other genotypes. Second, the genetic makeup of an individual is defined by many genes. Some phenotypes might depend on multiple genes. In other settings, we might be interested in multiple phenotypes, which (naturally) implies a dependence on several genes. Finally, as we now discuss, the inheritance patterns of different genes are not independent of each other.

Recall that each of the person's autosomal chromosome is inherited from one of her parents. However, each of the parents also has two copies of each autosomal chromosome. These two copies, within each parent, recombine to produce the chromosome that is transmitted to the child. Thus, the maternal chromosome inherited by Bart is a combination of the chromosomes inherited by his mother Marge from her mother Jackie and her father Clancy. The recombination process is stochastic, but only a handful recombination events take place within a chromosome in a single generation. Thus, if Bart inherited the allele for some locus from the chromosome his mother inherited from her mother Jackie, he is also much more likely to inherit Jackie's copy for a nearby locus. Thus, to construct an appropriate model for multilocus inheritance, we must take into consideration the probability of a recombination taking place between pairs of adjacent loci.

We can facilitate this modeling by introducing selector variables that capture the inheritance pattern along the chromosome. In particular, for each locus  $\ell$  and each child  $c$ , we have a variable  $S(\ell, c, m)$  that takes the value 1 if the locus  $\ell$  in  $c$ 's maternal chromosome was inherited from  $c$ 's maternal grandmother, and 2 if this locus was inherited from  $c$ 's maternal grandfather. We have a similar selector variable  $S(\ell, c, p)$  for  $c$ 's paternal chromosome. We can now model correlations induced by low recombination frequency by correlating the variables  $S(\ell, c, m)$  and  $S(\ell', c, m)$  for adjacent loci  $\ell, \ell'$ .

This type of model has been used extensively for many applications. In genetic counseling and prediction, one takes a phenotype with known loci and a set of observed phenotype and genotype data for some individuals in the pedigree to infer the genotype and phenotype for another person in the pedigree (say, a planned child). The genetic data can consist of direct measurements of the relevant disease loci (for some individuals) or measurements of nearby loci, which are correlated with the disease loci.

In linkage analysis, the task is a harder one: identifying the location of disease genes from pedigree data using some number of pedigrees where a large fraction of the individuals exhibit a disease phenotype. Here, the available data includes phenotype information for many individuals in the pedigree, as well as genotype information for loci whose location in the chromosome is known. Using the inheritance model, the researchers can evaluate the likelihood of these observations under different hypotheses about the location of the disease gene relative to the known loci. By repeated calculation of the probabilities in the network for different hypotheses, researchers can pinpoint the area that is "linked" to the disease. This much smaller region can then be used as the starting point for more detailed examination of genes in that area. This process is crucial, for it can allow the researchers to focus on a small area (for example, 1/10,000 of the genome).

As we will see in later chapters, the ability to describe the genetic inheritance process using a sparse Bayesian network provides us the capability to use sophisticated inference algorithms that allow us to reason about large pedigrees and multiple loci. It also allows us to use algorithms for model learning to obtain a deeper understanding of the genetic inheritance process, such as recombination rates in different regions or penetrance probabilities for different diseases.

### 3.2.3 Graphs and Distributions

The formal semantics of a Bayesian network graph is as a set of independence assertions. On the other hand, our Student BN was a graph annotated with CPDs, which defined a joint distribution via the chain rule for Bayesian networks. In this section, we show that these two definitions are, in fact, equivalent. A distribution  $P$  satisfies the local independencies associated with a graph  $\mathcal{G}$  if and only if  $P$  is representable as a set of CPDs associated with the graph  $\mathcal{G}$ . We begin by formalizing the basic concepts.

#### 3.2.3.1 I-Maps

We first define the set of independencies associated with a distribution  $P$ .

**Definition 3.2**  
independencies  
in  $P$

Let  $P$  be a distribution over  $\mathcal{X}$ . We define  $\mathcal{I}(P)$  to be the set of independence assertions of the form  $(X \perp Y \mid Z)$  that hold in  $P$ . ■

We can now rewrite the statement that " $P$  satisfies the local independencies associated with  $\mathcal{G}$ " simply as  $\mathcal{I}_\ell(\mathcal{G}) \subseteq \mathcal{I}(P)$ . In this case, we say that  $\mathcal{G}$  is an *I-map* (independency map) for  $P$ . However, it is useful to define this concept more broadly, since different variants of it will be used throughout the book.

**Definition 3.3**  
I-map

Let  $\mathcal{K}$  be any graph object associated with a set of independencies  $\mathcal{I}(\mathcal{K})$ . We say that  $\mathcal{K}$  is an I-map for a set of independencies  $\mathcal{I}$  if  $\mathcal{I}(\mathcal{K}) \subseteq \mathcal{I}$ . ■



We now say that  $\mathcal{G}$  is an I-map for  $P$  if  $\mathcal{G}$  is an I-map for  $\mathcal{I}(P)$ .

As we can see from the direction of the inclusion, for  $\mathcal{G}$  to be an I-map of  $P$ , it is necessary that  $\mathcal{G}$  does not mislead us regarding independencies in  $P$ : any independence that  $\mathcal{G}$  asserts must also hold in  $P$ . Conversely,  $P$  may have additional independencies that are not reflected in  $\mathcal{G}$ .

Let us illustrate the concept of an I-map on a very simple example.

### Example 3.1

Consider a joint probability space over two independent random variables  $X$  and  $Y$ . There are three possible graphs over these two nodes:  $\mathcal{G}_\emptyset$ , which is a disconnected pair  $X \quad Y$ ;  $\mathcal{G}_{X \rightarrow Y}$ , which has the edge  $X \rightarrow Y$ ; and  $\mathcal{G}_{Y \rightarrow X}$ , which contains  $Y \rightarrow X$ . The graph  $\mathcal{G}_\emptyset$  encodes the assumption that  $(X \perp Y)$ . The latter two encode no independence assumptions.

Consider the following two distributions:

$X$	$Y$	$P(X, Y)$
$x^0$	$y^0$	0.08
$x^0$	$y^1$	0.32
$x^1$	$y^0$	0.12
$x^1$	$y^1$	0.48

$X$	$Y$	$P(X, Y)$
$x^0$	$y^0$	0.4
$x^0$	$y^1$	0.3
$x^1$	$y^0$	0.2
$x^1$	$y^1$	0.1

In the example on the left,  $X$  and  $Y$  are independent in  $P$ ; for example,  $P(x^1) = 0.48 + 0.12 = 0.6$ ,  $P(y^1) = 0.8$ , and  $P(x^1, y^1) = 0.48 = 0.6 \cdot 0.8$ . Thus,  $(X \perp Y) \in \mathcal{I}(P)$ , and we have that  $\mathcal{G}_\emptyset$  is an I-map of  $P$ . In fact, all three graphs are I-maps of  $P$ :  $\mathcal{I}_\ell(\mathcal{G}_{X \rightarrow Y})$  is empty, so that trivially  $P$  satisfies all the independencies in it (similarly for  $\mathcal{G}_{Y \rightarrow X}$ ). In the example on the right,  $(X \perp Y) \notin \mathcal{I}(P)$ , so that  $\mathcal{G}_\emptyset$  is not an I-map of  $P$ . Both other graphs are I-maps of  $P$ . ■

### 3.2.3.2 I-Map to Factorization

A BN structure  $\mathcal{G}$  encodes a set of conditional independence assumptions; every distribution for which  $\mathcal{G}$  is an I-map must satisfy these assumptions. This property is the key to allowing the compact factorized representation that we saw in the Student example in section 3.2.1. The basic principle is the same as the one we used in the naive Bayes decomposition in section 3.1.3.

Consider any distribution  $P$  for which our Student BN  $G_{student}$  is an I-map. We will decompose the joint and show that it factorizes into local probabilistic models, as in section 3.2.1. Consider the joint distribution  $P(I, D, G, L, S)$ ; from the chain rule for probabilities (equation (2.5)). We can decompose this joint in the following way:

$$P(I, D, G, L, S) = P(I)P(D | I)P(G | I, D)P(L | I, D, G)P(S | I, D, G, L). \quad (3.15)$$

This transformation relies on no assumptions; it holds for any joint distribution  $P$ . However, it is also not very helpful, since the conditional probabilities in the factorization on the right-hand side are neither natural nor compact. For example, the last factor requires the specification of 24 conditional probabilities:  $P(s^1 | i, d, g, l)$  for every assignment of values  $i, d, g, l$ .

This form, however, allows us to apply the conditional independence assumptions induced from the BN. Let us assume that  $G_{student}$  is an I-map for our distribution  $P$ . In particular, from equation (3.14), we have that  $(I \perp D) \in \mathcal{I}(P)$ . From that, we can conclude that  $P(D | I) = P(D)$ , allowing us to simplify the second factor on the right-hand side. Similarly, we know from

equation (3.10) that  $(L \perp I, D \mid G) \in \mathcal{I}(P)$ . Hence,  $P(L \mid I, D, G) = P(L \mid G)$ , allowing us to simplify the third term. Using equation (3.11) in a similar way, we obtain that

$$P(I, D, G, L, S) = P(I)P(D)P(G \mid I, D)P(L \mid G)P(S \mid I). \quad (3.16)$$

This factorization is precisely the one we used in section 3.2.1.

This result tells us that any entry in the joint distribution can be computed as a product of factors, one for each variable. Each factor represents a conditional probability of the variable given its parents in the network. This factorization applies to *any* distribution  $P$  for which  $G_{student}$  is an I-map.

We now state and prove this fundamental result more formally.

**Definition 3.4**  
factorization

Let  $\mathcal{G}$  be a BN graph over the variables  $X_1, \dots, X_n$ . We say that a distribution  $P$  over the same space factorizes according to  $\mathcal{G}$  if  $P$  can be expressed as a product

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid \text{Pa}_{X_i}^{\mathcal{G}}). \quad (3.17)$$

chain rule for  
Bayesian  
networks

This equation is called the chain rule for Bayesian networks. The individual factors  $P(X_i \mid \text{Pa}_{X_i}^{\mathcal{G}})$  are called conditional probability distributions (CPDs) or local probabilistic models. ■

CPD

**Definition 3.5**

Bayesian network

A Bayesian network is a pair  $B = (\mathcal{G}, P)$  where  $P$  factorizes over  $\mathcal{G}$ , and where  $P$  is specified as a set of CPDs associated with  $\mathcal{G}$ 's nodes. The distribution  $P$  is often annotated  $P_B$ . ■

We can now prove that the phenomenon we observed for  $G_{student}$  holds more generally.

**Theorem 3.1**

Let  $\mathcal{G}$  be a BN structure over a set of random variables  $\mathcal{X}$ , and let  $P$  be a joint distribution over the same space. If  $\mathcal{G}$  is an I-map for  $P$ , then  $P$  factorizes according to  $\mathcal{G}$ .

topological  
ordering

**PROOF** Assume, without loss of generality, that  $X_1, \dots, X_n$  is a topological ordering of the variables in  $\mathcal{X}$  relative to  $\mathcal{G}$  (see definition 2.19). As in our example, we first use the chain rule for probabilities:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid X_1, \dots, X_{i-1}).$$

Now, consider one of the factors  $P(X_i \mid X_1, \dots, X_{i-1})$ . As  $\mathcal{G}$  is an I-map for  $P$ , we have that  $(X_i \perp \text{NonDescendants}_{X_i} \mid \text{Pa}_{X_i}^{\mathcal{G}}) \in \mathcal{I}(P)$ . By assumption, all of  $X_i$ 's parents are in the set  $X_1, \dots, X_{i-1}$ . Furthermore, none of  $X_i$ 's descendants can possibly be in the set. Hence,

$$\{X_1, \dots, X_{i-1}\} = \text{Pa}_{X_i} \cup Z$$

where  $Z \subseteq \text{NonDescendants}_{X_i}$ . From the local independencies for  $X_i$  and from exercise 2.9(1) it follows that  $(X_i \perp Z \mid \text{Pa}_{X_i})$ . Hence, we have that

$$P(X_i \mid X_1, \dots, X_{i-1}) = P(X_i \mid \text{Pa}_{X_i}).$$

Applying this transformation to all of the factors in the chain rule decomposition, the result follows.

Thus, the conditional independence assumptions implied by a BN structure  $\mathcal{G}$  allow us to factorize a distribution  $P$  for which  $\mathcal{G}$  is an I-map into small CPDs. Note that the proof is constructive, providing a precise algorithm for constructing the factorization given the distribution  $P$  and the graph  $\mathcal{G}$ .

The resulting factorized representation can be substantially more compact, particularly for sparse structures.

**Example 3.2**

*In our Student example, the number of independent parameters is fifteen: we have two binomial distributions  $P(I)$  and  $P(D)$ , with one independent parameter each; we have four multinomial distributions over  $\mathcal{G}$  — one for each assignment of values to  $I$  and  $D$  — each with two independent parameters; we have three binomial distributions over  $L$ , each with one independent parameter; and similarly two binomial distributions over  $S$ , each with an independent parameter. The specification of the full joint distribution would require  $48 - 1 = 47$  independent parameters. ■*

More generally, in a distribution over  $n$  binary random variables, the specification of the joint distribution requires  $2^n - 1$  independent parameters. If the distribution factorizes according to a graph  $\mathcal{G}$  where each node has at most  $k$  parents, the total number of independent parameters required is less than  $n \cdot 2^k$  (see exercise 3.6). In many applications, we can assume a certain locality of influence between variables: although each variable is generally correlated with many of the others, it often depends *directly* on only a small number of other variables. Thus, in many cases,  $k$  will be very small, even though  $n$  is large. As a consequence, the number of parameters in the Bayesian network representation is typically exponentially smaller than the number of parameters of a joint distribution. This property is one of the main benefits of the Bayesian network representation.

**3.2.3.3 Factorization to I-Map**

Theorem 3.1 shows one direction of the fundamental connection between the conditional independencies encoded by the BN structure and the factorization of the distribution into local probability models: that the conditional independencies imply factorization. The converse also holds: factorization according to  $\mathcal{G}$  implies the associated conditional independencies.

**Theorem 3.2**

*Let  $\mathcal{G}$  be a BN structure over a set of random variables  $\mathcal{X}$  and let  $P$  be a joint distribution over the same space. If  $P$  factorizes according to  $\mathcal{G}$ , then  $\mathcal{G}$  is an I-map for  $P$ .*

We illustrate this theorem by example, leaving the proof as an exercise (exercise 3.9). Let  $P$  be some distribution that factorizes according to  $G_{\text{student}}$ . We need to show that  $\mathcal{I}_\ell(G_{\text{student}})$  holds in  $P$ . Consider the independence assumption for the random variable  $S$  —  $(S \perp D, G, L \mid I)$ . To prove that it holds for  $P$ , we need to show that

$$P(S \mid I, D, G, L) = P(S \mid I).$$

By definition,

$$P(S \mid I, D, G, L) = \frac{P(S, I, D, G, L)}{P(I, D, G, L)}.$$

By the chain rule for BNs equation (3.16), the numerator is equal to  $P(I)P(D)P(G | I, D)P(L | G)P(S | I)$ . By the process of marginalizing over a joint distribution, we have that the denominator is:

$$\begin{aligned} P(I, D, G, L) &= \sum_S P(I, D, G, L, S) \\ &= \sum_S P(I)P(D)P(G | I, D)P(L | G)P(S | I) \\ &= P(I)P(D)P(G | I, D)P(L | G) \sum_S P(S | I) \\ &= P(I)P(D)P(G | I, D)P(L | G), \end{aligned}$$

where the last step is a consequence of the fact that  $P(S | I)$  is a distribution over values of  $S$ , and therefore it sums to 1. We therefore have that

$$\begin{aligned} P(S | I, D, G, L) &= \frac{P(S, I, D, G, L)}{P(I, D, G, L)} \\ &= \frac{P(I)P(D)P(G | I, D)P(L | G)P(S | I)}{P(I)P(D)P(G | I, D)P(L | G)} \\ &= P(S | I). \end{aligned}$$

---

**Box 3.C — Skill: Knowledge Engineering.** *Our discussion of Bayesian network construction focuses on the process of going from a given distribution to a Bayesian network. Real life is not like that. We have a vague model of the world, and we need to crystallize it into a network structure and parameters. This task breaks down into several components, each of which can be quite subtle. Unfortunately, modeling mistakes can have significant consequences for the quality of the answers obtained from the network, or to the cost of using the network in practice.*

**Picking variables** When we model a domain, there are many possible ways to describe the relevant entities and their attributes. Choosing which random variables to use in the model is often one of the hardest tasks, and this decision has implications throughout the model. A common problem is using ill-defined variables. For example, deciding to include the variable *Fever* to describe a patient in a medical domain seems fairly innocuous. However, does this random variable relate to the internal temperature of the patient? To the thermometer reading (if one is taken by the medical staff)? Does it refer to the temperature of the patient at a specific moment (for example, the time of admission to the hospital) or to occurrence of a fever over a prolonged period? Clearly, each of these might be a reasonable attribute to model, but the interaction of *Fever* with other variables depends on the specific interpretation we use.

clarity test

As this example shows, we must be precise in defining the variables in the model. The clarity test is a good way of evaluating whether they are sufficiently well defined. Assume that we are a million years after the events described in the domain; can an omniscient being, one who saw everything, determine the value of the variable? For example, consider a *Weather* variable with a value *sunny*. To be absolutely precise, we must define where we check the weather, at what time,

and what fraction of the sky must be clear in order for it to be sunny. For a variable such as Heart-attack, we must specify how large the heart attack has to be, during what period of time it has to happen, and so on. By contrast, a variable such as Risk-of-heart-attack is meaningless, as even an omniscient being cannot evaluate whether a person had high risk or low risk, only whether the heart attack occurred or not. Introducing variables such as this confounds actual events and their probability. Note, however, that we can use a notion of "risk group," as long as it is defined in terms of clearly specified attributes such as age or lifestyle.

If we are not careful in our choice of variables, we will have a hard time making sure that evidence observed and conclusions made are coherent.

Generally speaking, we want our model to contain variables that we can potentially observe or that we may want to query. However, sometimes we want to put in a hidden variable that is neither observed nor directly of interest. Why would we want to do that? Let us go back to the example of the cholesterol test. For the answers to be accurate, the subject has to have eaten nothing after 10:00 PM the previous evening. If the person eats (having no willpower), the results are consistently off. We do not really care about a Willpower variable, nor can we observe it. However, without it, all of the different cholesterol tests become correlated. To avoid graphs where all the tests are correlated, it is better to put in this additional hidden variable, rendering them conditionally independent given the true cholesterol level and the person's willpower.

On the other hand, it is not necessary to add every variable that might be relevant. In our Student example, the student's SAT score may be affected by whether he goes out for drinks on the night before the exam. Is this variable important to represent? The probabilities already account for the fact that he may achieve a poor score despite being intelligent. It might not be worthwhile to include this variable if it cannot be observed.

It is also important to specify a reasonable domain of values for our variables. In particular, if our partition is not fine enough, conditional independence assumptions may be false. For example, we might want to construct a model where we have a person's cholesterol level, and two cholesterol tests that are conditionally independent given the person's true cholesterol level. We might choose to define the value normal to correspond to levels up to 200, and high to levels above 200. But it may be the case that both tests are more likely to fail if the person's cholesterol is marginal (200-240). In this case, the assumption of conditional independence given the value (high/normal) of the cholesterol test is false. It is only true if we add a marginal value.

**Picking structure** As we saw, there are many structures that are consistent with the same set of independencies. One successful approach is to choose a structure that reflects the causal order and dependencies, so that causes are parents of the effect. Such structures tend to work well. Either because of some real locality of influence in the world, or because of the way people perceive the world, causal graphs tend to be sparser. It is important to stress that the causality is in the world, not in our inference process. For example, in an automobile insurance network, it is tempting to put Previous-accident as a parent of Good-driver, because that is how the insurance company thinks about the problem. This is not the causal order in the world, because being a bad driver causes previous (and future) accidents. In principle, there is nothing to prevent us from directing the edges in this way. However, a noncausal ordering often requires that we introduce many additional edges to account for induced dependencies (see section 3.4.1).

One common approach to constructing a structure is a backward construction process. We begin with a variable of interest, say Lung-Cancer. We then try to elicit a prior probability for that

hidden variable

variable. If our expert responds that this probability is not determinable, because it depends on other factors, that is a good indication that these other factors should be added as parents for that variable (and as variables into the network). For example, we might conclude using this process that Lung-Cancer really should have Smoking as a parent, and (perhaps not as obvious) that Smoking should have Gender and Age as a parent. This approach, called extending the conversation, avoids probability estimates that result from an average over a heterogeneous population, and therefore leads to more precise probability estimates.

When determining the structure, however, we must also keep in mind that approximations are inevitable. For many pairs of variables, we can construct a scenario where one depends on the other. For example, perhaps Difficulty depends on Intelligence, because the professor is more likely to make a class difficult if intelligent students are registered. In general, there are many weak influences that we might choose to model, but if we put in all of them, the network can become very complex. Such networks are problematic from a representational perspective: they are hard to understand and hard to debug, and eliciting (or learning) parameters can get very difficult. Moreover, as reasoning in Bayesian networks depends strongly on their connectivity (see section 9.4), adding such edges can make the network too expensive to use.

This final consideration may lead us, in fact, to make approximations that we know to be wrong. For example, in networks for fault or medical diagnosis, the correct approach is usually to model each possible fault as a separate random variable, allowing for multiple failures. However, such networks might be too complex to perform effective inference in certain settings, and so we may sometimes resort to a single fault approximation, where we have a single random variable encoding the primary fault or disease.

**Picking probabilities** One of the most challenging tasks in constructing a network manually is eliciting probabilities from people. This task is somewhat easier in the context of causal models, since the parameters tend to be natural and more interpretable. Nevertheless, people generally dislike committing to an exact estimate of probability.

One approach is to elicit estimates qualitatively, using abstract terms such as "common," "rare," and "surprising," and then assign these to numbers using a predefined scale. This approach is fairly crude, and often can lead to misinterpretation. There are several approaches developed for assisting in eliciting probabilities from people. For example, one can visualize the probability of the event as an area (slice of a pie), or ask people how they would compare the probability in question to certain predefined lotteries. Nevertheless, probability elicitation is a long, difficult process, and one whose outcomes are not always reliable: the elicitation method can often influence the results, and asking the same question using different phrasing can often lead to significant differences in the answer. For example, studies show that people's estimates for an event such as "Death by disease" are significantly lower than their estimates for this event when it is broken down into different possibilities such as "Death from cancer," "Death from heart disease," and so on.

How important is it that we get our probability estimates exactly right? In some cases, small errors have very little effect. For example, changing a conditional probability of 0.7 to 0.75 generally does not have a significant effect. Other errors, however, can have a significant effect:

- **Zero probabilities:** A common mistake is to assign a probability of zero to an event that is extremely unlikely, but not impossible. The problem is that one can never condition away a zero probability, no matter how much evidence we get. When an event is unlikely

but not impossible, giving it probability zero is guaranteed to lead to irrecoverable errors. For example, in one of the early versions of the the Pathfinder system (box 3.D), 10 percent of the misdiagnoses were due to zero probability estimates given by the expert to events that were unlikely but not impossible. As a general rule, very few things (except definitions) have probability zero, and we must be careful in assigning zeros.

- **Orders of magnitude:** Small differences in very low probability events can make a large difference to the network conclusions. Thus, a (conditional) probability of  $10^{-4}$  is very different from  $10^{-5}$ .
- **Relative values:** The qualitative behavior of the conclusions reached by the network — the value that has the highest probability — is fairly sensitive to the relative sizes of  $P(x | y)$  for different values  $y$  of  $\text{Pa}_x$ . For example, it is important that the network encode correctly that the probability of having a high fever is greater when the patient has pneumonia than when he has the flu.

A very useful tool for estimating network parameters is sensitivity analysis, which allows us to determine the extent to which a given probability parameter affects the outcome. This process allows us to evaluate whether it is important to get a particular CPD entry right. It also helps us figure out which CPD entries are responsible for an answer to some query that does not match our intuitions.

sensitivity  
analysis

medical diagnosis  
expert system

---

**Box 3.D — Case Study: Medical Diagnosis Systems.** One of the earliest applications of Bayesian networks was to the task of medical diagnosis. In the 1980s, a very active area of research was the construction of expert systems — computer-based systems that replace or assist an expert in performing a complex task. One such task that was tackled in several ways was medical diagnosis. This task, more than many others, required a treatment of uncertainty, due to the complex, nondeterministic relationships between findings and diseases. Thus, it formed the basis for experimentation with various formalisms for uncertain reasoning.

Pathfinder

The Pathfinder expert system was designed by Heckerman and colleagues (Heckerman and Nathwani 1992a; Heckerman et al. 1992; Heckerman and Nathwani 1992b) to help a pathologist diagnose diseases in lymph nodes. Ultimately, the model contained more than sixty different diseases and around a hundred different features. It evolved through several versions, including some based on nonprobabilistic formalisms, and several that used variants of Bayesian networks. Its diagnostic ability was evaluated over real pathological cases and compared to the diagnoses of pathological experts.

One of the first models used was a simple naive Bayes model, which was compared to the models based on alternative uncertainty formalisms, and judged to be superior in its diagnostic ability. It therefore formed the basis for subsequent development of the system.

The same evaluation pointed out important problems in the way in which parameters were elicited from the expert. First, it was shown that 10 percent of the cases were diagnosed incorrectly, because the correct disease was ruled out by a finding that was unlikely, but not impossible, to manifest in that disease. Second, in the original construction, the expert estimated the probabilities  $P(\text{Finding} | \text{Disease})$  by fixing a single disease and evaluating the probabilities of all its findings. It was found that the expert was more comfortable considering a single finding and evaluating its

probability across all diseases. This approach allows the expert to compare the relative values of the same finding across multiple diseases, as described in box 3.C.

With these two lessons in mind, another version of Pathfinder — Pathfinder III — was constructed, still using the naive Bayes model. Finally, Pathfinder IV used a full Bayesian network, with a single disease hypothesis but with dependencies between the features. Pathfinder IV was constructed using a similarity network (see box 5.B), significantly reducing the number of parameters that must be elicited. Pathfinder IV, viewed as a Bayesian network, had a total of around 75,000 parameters, but the use of similarity networks allowed the model to be constructed with fewer than 14,000 distinct parameters. Overall, the structure of Pathfinder IV took about 35 hours to define, and the parameters 40 hours.

A comprehensive evaluation of the performance of the two models revealed some important insights. First, the Bayesian network performed as well or better on most cases than the naive Bayes model. In most of the cases where the Bayesian network performed better, the use of richer dependency models was a contributing factor. As expected, these models were useful because they address the strong conditional independence assumptions of the naive Bayes model, as described in box 3.A. Somewhat more surprising, they also helped in allowing the expert to condition the probabilities on relevant factors other than the disease, using the process of extending the conversation described in box 3.C, leading to more accurate elicited probabilities. Finally, the use of similarity networks led to more accurate models, for the smaller number of elicited parameters reduced irrelevant fluctuations in parameter values (due to expert inconsistency) that can lead to spurious dependencies.

Overall, the Bayesian network model agreed with the predictions of an expert pathologist in 50/53 cases, as compared with 47/53 cases for the naive Bayes model, with significant therapeutic implications. A later evaluation showed that the diagnostic accuracy of Pathfinder IV was at least as good as that of the expert used to design the system. When used with less expert pathologists, the system significantly improved the diagnostic accuracy of the physicians alone. Moreover, the system showed greater ability to identify important findings and to integrate these findings into a correct diagnosis.

Unfortunately, multiple reasons prevent the widespread adoption of Bayesian networks as an aid for medical diagnosis, including legal liability issues for misdiagnoses and incompatibility with the physicians' workflow. However, several such systems have been fielded, with significant success. Moreover, similar technology is being used successfully in a variety of other diagnosis applications (see box 23.C).

---

### 3.3 Independencies in Graphs

Dependencies and independencies are key properties of a distribution and are crucial for understanding its behavior. As we will see, independence properties are also important for answering queries: they can be exploited to reduce substantially the computation cost of inference. Therefore, it is important that our representations make these properties clearly visible both to a user and to algorithms that manipulate the BN data structure.

As we discussed, a graph structure  $\mathcal{G}$  encodes a certain set of conditional independence assumptions  $\mathcal{I}_\ell(\mathcal{G})$ . Knowing only that a distribution  $P$  factorizes over  $\mathcal{G}$ , we can conclude



that it satisfies  $\mathcal{I}_\ell(\mathcal{G})$ . An immediate question is whether there are other independencies that we can “read off” directly from  $\mathcal{G}$ . That is, are there other independencies that hold for *every* distribution  $P$  that factorizes over  $\mathcal{G}$ ?

### 3.3.1 D-separation

Our aim in this section is to understand when we can *guarantee* that an independence  $(X \perp Y \mid Z)$  holds in a distribution associated with a BN structure  $\mathcal{G}$ . To understand when a property is guaranteed to hold, it helps to consider its converse: “Can we imagine a case where it does not?” Thus, we focus our discussion on analyzing when it is *possible* that  $X$  can influence  $Y$  given  $Z$ . If we construct an example where this influence occurs, then the converse property  $(X \perp Y \mid Z)$  cannot hold for all of the distributions that factorize over  $\mathcal{G}$ , and hence the independence property  $(X \perp Y \mid Z)$  cannot follow from  $\mathcal{I}_\ell(\mathcal{G})$ .

We therefore begin with an intuitive case analysis: Here, we try to understand when an observation regarding a variable  $X$  can possibly change our beliefs about  $Y$ , in the presence of evidence about the variables  $Z$ . Although this analysis will be purely intuitive, we will show later that our conclusions are actually provably correct.

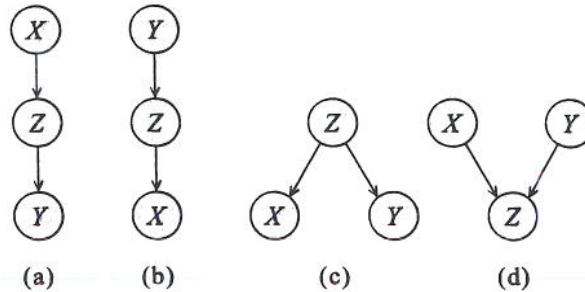
**Direct connection** We begin with the simple case, when  $X$  and  $Y$  are directly connected via an edge, say  $X \rightarrow Y$ . For any network structure  $\mathcal{G}$  that contains the edge  $X \rightarrow Y$ , it is possible to construct a distribution where  $X$  and  $Y$  are correlated regardless of any evidence about any of the other variables in the network. In other words, if  $X$  and  $Y$  are directly connected, we can always get examples where they influence each other, regardless of  $Z$ .

In particular, assume that  $Val(X) = Val(Y)$ ; we can simply set  $X = Y$ . That, by itself, however, is not enough; if (given the evidence  $Z$ )  $X$  deterministically takes some particular value, say 0, then  $X$  and  $Y$  both deterministically take that value, and are uncorrelated. We therefore set the network so that  $X$  is (for example) uniformly distributed, regardless of the values of any of its parents. This construction suffices to induce a correlation between  $X$  and  $Y$ , regardless of the evidence.

**Indirect connection** Now consider the more complicated case when  $X$  and  $Y$  are not directly connected, but there is a trail between them in the graph. We begin by considering the simplest such case: a three-node network, where  $X$  and  $Y$  are not directly connected, but where there is a trail between them via  $Z$ . It turns out that this simple case is the key to understanding the whole notion of indirect interaction in Bayesian networks.

There are four cases where  $X$  and  $Y$  are connected via  $Z$ , as shown in figure 3.5. The first two correspond to causal chains (in either direction), the third to a common cause, and the fourth to a common effect. We analyze each in turn.

**Indirect causal effect (figure 3.5a).** To gain intuition, let us return to the Student example, where we had a causal trail  $I \rightarrow G \rightarrow L$ . Let us begin with the case where  $G$  is not observed. Intuitively, if we observe that the student is intelligent, we are more inclined to believe that he gets an A, and therefore that his recommendation letter is strong. In other words, the probability of these latter events is higher conditioned on the observation that the student is intelligent. In fact, we saw precisely this behavior in the distribution of figure 3.4. Thus, in this case, we believe that  $X$  can influence  $Y$  via  $Z$ .



**Figure 3.5** The four possible two-edge trails from  $X$  to  $Y$  via  $Z$ : (a) An indirect causal effect; (b) An indirect evidential effect; (c) A common cause; (d) A common effect.

Now assume that  $Z$  is observed, that is,  $Z \in \mathcal{Z}$ . As we saw in our analysis of the Student example, if we observe the student's grade, then (as we assumed) his intelligence no longer influences his letter. In fact, the local independencies for this network tell us that  $(L \perp I \mid G)$ . Thus, we conclude that  $X$  cannot influence  $Y$  via  $Z$  if  $Z$  is observed.

**Indirect evidential effect (figure 3.5b).** Returning to the Student example, we have a chain  $I \rightarrow G \rightarrow L$ . We have already seen that observing a strong recommendation letter for the student changes our beliefs in his intelligence. Conversely, once the grade is observed, the letter gives no additional information about the student's intelligence. Thus, our analysis in the case  $Y \rightarrow Z \rightarrow X$  here is identical to the causal case:  $X$  can influence  $Y$  via  $Z$ , but only if  $Z$  is not observed. The similarity is not surprising, as dependence is a symmetrical notion. Specifically, if  $(X \perp Y)$  does not hold, then  $(Y \perp X)$  does not hold either.

**Common cause (figure 3.5c).** This case is one that we have analyzed extensively, both within the simple naive Bayes model of section 3.1.3 and within our Student example. Our example has the student's intelligence  $I$  as a parent of his grade  $G$  and his SAT score  $S$ . As we discussed,  $S$  and  $G$  are correlated in this model, in that observing (say) a high SAT score gives us information about a student's intelligence and hence helps us predict his grade. However, once we observe  $I$ , this correlation disappears, and  $S$  gives us no additional information about  $G$ . Once again, for this network, this conclusion follows from the local independence assumption for the node  $G$  (or for  $S$ ). Thus, our conclusion here is identical to the previous two cases:  $X$  can influence  $Y$  via  $Z$  if and only if  $Z$  is not observed.

**Common effect (figure 3.5d).** In all of the three previous cases, we have seen a common pattern:  $X$  can influence  $Y$  via  $Z$  if and only if  $Z$  is not observed. Therefore, we might expect that this pattern is universal, and will continue through this last case. Somewhat surprisingly, this is not the case. Let us return to the Student example and consider  $I$  and  $D$ , which are parents of  $G$ . When  $G$  is not observed, we have that  $I$  and  $D$  are independent. In fact, this conclusion follows (once again) from the local independencies from the network. Thus, in this case, influence cannot "flow" along the trail  $X \rightarrow Z \leftarrow Y$  if the intermediate node  $Z$  is not observed.

On the other hand, consider the behavior when  $Z$  is observed. In our discussion of the Student example, we analyzed precisely this case, which we called intercausal reasoning; we showed, for example, that the probability that the student has high intelligence goes down

dramatically when we observe that his grade is a C ( $G = g^3$ ), but then goes up when we observe that the class is a difficult one  $D = d^1$ . Thus, in presence of the evidence  $G = g^3$ , we have that  $I$  and  $D$  are correlated.

Let us consider a variant of this last case. Assume that we do not observe the student's grade, but we do observe that he received a weak recommendation letter ( $L = l^0$ ). Intuitively, the same phenomenon happens. The weak letter is an indicator that he received a low grade, and therefore it suffices to correlate  $I$  and  $D$ .

When influence can flow from  $X$  to  $Y$  via  $Z$ , we say that the trail  $X \rightleftharpoons Z \rightleftharpoons Y$  is *active*. The results of our analysis for active two-edge trails are summarized thus:

- **Causal trail**  $X \rightarrow Z \rightarrow Y$ : active if and only if  $Z$  is not observed.
- **Evidential trail**  $X \leftarrow Z \leftarrow Y$ : active if and only if  $Z$  is not observed.
- **Common cause**  $X \leftarrow Z \rightarrow Y$ : active if and only if  $Z$  is not observed.
- **Common effect**  $X \rightarrow Z \leftarrow Y$ : active if and only if either  $Z$  or one of  $Z$ 's descendants is observed.

v-structure

A structure where  $X \rightarrow Z \leftarrow Y$  (as in figure 3.5d) is also called a *v-structure*.

It is useful to view probabilistic influence as a flow in the graph. Our analysis here tells us when influence from  $X$  can "flow" through  $Z$  to affect our beliefs about  $Y$ .

**General Case** Now consider the case of a longer trail  $X_1 \rightleftharpoons \dots \rightleftharpoons X_n$ . Intuitively, for influence to "flow" from  $X_1$  to  $X_n$ , it needs to flow through every single node on the trail. In other words,  $X_1$  can influence  $X_n$  if every two-edge trail  $X_{i-1} \rightleftharpoons X_i \rightleftharpoons X_{i+1}$  along the trail allows influence to flow.

We can summarize this intuition in the following definition:

Definition 3.6

observed variable  
active trail

Let  $\mathcal{G}$  be a BN structure, and  $X_1 \rightleftharpoons \dots \rightleftharpoons X_n$  a trail in  $\mathcal{G}$ . Let  $Z$  be a subset of observed variables. The trail  $X_1 \rightleftharpoons \dots \rightleftharpoons X_n$  is active given  $Z$  if

- Whenever we have a v-structure  $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$ , then  $X_i$  or one of its descendants are in  $Z$ ;
- no other node along the trail is in  $Z$ .

Note that if  $X_1$  or  $X_n$  are in  $Z$  the trail is not active. ■

In our Student BN, we have that  $D \rightarrow G \leftarrow I \rightarrow S$  is not an active trail for  $Z = \emptyset$ , because the v-structure  $D \rightarrow G \leftarrow I$  is not activated. That same trail is active when  $Z = \{L\}$ , because observing the descendant of  $G$  activates the v-structure. On the other hand, when  $Z = \{L, I\}$ , the trail is not active, because observing  $I$  blocks the trail  $G \leftarrow I \rightarrow S$ .

What about graphs where there is more than one trail between two nodes? Our flow intuition continues to carry through: one node can influence another if there is any trail along which influence can flow. Putting these intuitions together, we obtain the notion of *d-separation*, which provides us with a notion of separation between nodes in a directed graph (hence the term d-separation, for directed separation):

d-separation

Definition 3.7

Let  $X, Y, Z$  be three sets of nodes in  $\mathcal{G}$ . We say that  $X$  and  $Y$  are d-separated given  $Z$ , denoted  $d\text{-sep}_{\mathcal{G}}(X; Y \mid Z)$ , if there is no active trail between any node  $X \in X$  and  $Y \in Y$  given  $Z$ .

We use  $\mathcal{I}(\mathcal{G})$  to denote the set of independencies that correspond to d-separation:

$$\mathcal{I}(\mathcal{G}) = \{(X \perp Y \mid Z) : \text{d-sep}_{\mathcal{G}}(X; Y \mid Z)\}. \quad \blacksquare$$

global Markov  
independencies

This set is also called the set of *global Markov independencies*. The similarity between the notation  $\mathcal{I}(\mathcal{G})$  and our notation  $\mathcal{I}(P)$  is not coincidental: As we discuss later, the independencies in  $\mathcal{I}(\mathcal{G})$  are precisely those that are guaranteed to hold for every distribution over  $\mathcal{G}$ .

### 3.3.2 Soundness and Completeness

So far, our definition of d-separation has been based on our intuitions regarding flow of influence, and on our one example. As yet, we have no guarantee that this analysis is “correct.” Perhaps there is a distribution over the BN where  $X$  can influence  $Y$  despite the fact that all trails between them are blocked.

soundness of  
d-separation

Hence, the first property we want to ensure for d-separation as a method for determining independence is *soundness*: if we find that two nodes  $X$  and  $Y$  are d-separated given some  $Z$ , then we are guaranteed that they are, in fact, conditionally independent given  $Z$ .

Theorem 3.3

If a distribution  $P$  factorizes according to  $\mathcal{G}$ , then  $\mathcal{I}(\mathcal{G}) \subseteq \mathcal{I}(P)$ .

In other words, any independence reported by d-separation is satisfied by the underlying distribution. The proof of this theorem requires some additional machinery that we introduce in chapter 4, so we defer the proof to that chapter (see section 4.5.1.1).

completeness of  
d-separation

A second desirable property is the complementary one — *completeness*: d-separation detects *all* possible independencies. More precisely, if we have that two variables  $X$  and  $Y$  are independent given  $Z$ , then they are d-separated. A careful examination of the completeness property reveals that it is ill defined, inasmuch as it does not specify the distribution in which  $X$  and  $Y$  are independent.

To formalize this property, we first define the following notion:

Definition 3.8

faithful

A distribution  $P$  is faithful to  $\mathcal{G}$  if, whenever  $(X \perp Y \mid Z) \in \mathcal{I}(P)$ , then  $\text{d-sep}_{\mathcal{G}}(X; Y \mid Z)$ . In other words, any independence in  $P$  is reflected in the d-separation properties of the graph.  $\blacksquare$

We can now provide one candidate formalization of the completeness property is as follows:

- For any distribution  $P$  that factorizes over  $\mathcal{G}$ , we have that  $P$  is faithful to  $\mathcal{G}$ ; that is, if  $X$  and  $Y$  are *not* d-separated given  $Z$  in  $\mathcal{G}$ , then  $X$  and  $Y$  are *dependent* in *all* distributions  $P$  that factorize over  $\mathcal{G}$ .

This property is the obvious converse to our notion of soundness: If true, the two together would imply that, for any  $P$  that factorizes over  $\mathcal{G}$ , we have that  $\mathcal{I}(P) = \mathcal{I}(\mathcal{G})$ . Unfortunately, this highly desirable property is easily shown to be false: Even if a distribution factorizes over  $\mathcal{G}$ , it can still contain additional independencies that are not reflected in the structure.

Example 3.3

Consider a distribution  $P$  over two variables  $A$  and  $B$ , where  $A$  and  $B$  are independent. One possible I-map for  $P$  is the network  $A \rightarrow B$ . For example, we can set the CPD for  $B$  to be

	$b^0$	$b^1$
$a^0$	0.4	0.6
$a^1$	0.4	0.6

This example clearly violates the first candidate definition of completeness, because the graph  $\mathcal{G}$  is an I-map for the distribution  $P$ , yet there are independencies that hold for this distribution but do not follow from  $d$ -separation. In fact, these are not independencies that we can hope to discover by examining the network structure. ■

Thus, the completeness property does not hold for this candidate definition of completeness. We therefore adopt a weaker yet still useful definition:

- If  $(X \perp Y \mid Z)$  in all distributions  $P$  that factorize over  $\mathcal{G}$ , then  $d\text{-sep}_{\mathcal{G}}(X; Y \mid Z)$ . And the contrapositive: If  $X$  and  $Y$  are not  $d$ -separated given  $Z$  in  $\mathcal{G}$ , then  $X$  and  $Y$  are dependent in some distribution  $P$  that factorizes over  $\mathcal{G}$ .

Using this definition, we can show:

**Theorem 3.4**

Let  $\mathcal{G}$  be a BN structure. If  $X$  and  $Y$  are not  $d$ -separated given  $Z$  in  $\mathcal{G}$ , then  $X$  and  $Y$  are dependent given  $Z$  in some distribution  $P$  that factorizes over  $\mathcal{G}$ .

**PROOF** The proof constructs a distribution  $P$  that makes  $X$  and  $Y$  correlated. The construction is roughly as follows. As  $X$  and  $Y$  are not  $d$ -separated, there exists an active trail  $U_1, \dots, U_k$  between them. We define CPDs for the variables on the trail so as to make each pair  $U_i, U_{i+1}$  correlated; in the case of a  $v$ -structure  $U_i \rightarrow U_{i+1} \leftarrow U_{i+2}$ , we define the CPD of  $U_{i+1}$  so as to ensure correlation, and also define the CPDs of the path to some downstream evidence node, in a way that guarantees that the downstream evidence activates the correlation between  $U_i$  and  $U_{i+2}$ . All other CPDs in the graph are chosen to be uniform, and thus the construction guarantees that influence only flows along this single path, preventing cases where the influence of two (or more) paths cancel out. The details of the construction are quite technical and laborious, and we omit them. ■

We can view the completeness result as telling us that our definition of  $\mathcal{I}(\mathcal{G})$  is the maximal one. For any independence assertion that is not a consequence of  $d$ -separation in  $\mathcal{G}$ , we can always find a counterexample distribution  $P$  that factorizes over  $\mathcal{G}$ . In fact, this result can be strengthened significantly:

**Theorem 3.5**

For almost all distributions  $P$  that factorize over  $\mathcal{G}$ , that is, for all distributions except for a set of measure zero in the space of CPD parameterizations, we have that  $\mathcal{I}(P) = \mathcal{I}(\mathcal{G})$ .<sup>1</sup>

This result strengthens theorem 3.4 in two distinct ways: First, whereas theorem 3.4 shows that any dependency in the graph can be found in some distribution, this new result shows that there exists a single distribution that is faithful to the graph, that is, where all of the dependencies in the graph hold simultaneously. Second, not only does this property hold for a single distribution, but it also holds for almost all distributions that factorize over  $\mathcal{G}$ .

**PROOF** At a high level, the proof is based on the following argument: Each conditional independence assertion is a set of polynomial equalities over the space of CPD parameters (see

1. A set has measure zero if it is infinitesimally small relative to the overall space. For example, the set of all rationals has measure zero within the interval  $[0, 1]$ . A straight line has measure zero in the plane. This intuition is defined formally in the field of *measure theory*.

exercise 3.13). A basic property of polynomials is that a polynomial is either identically zero or it is nonzero almost everywhere (its set of roots has measure zero). Theorem 3.4 implies that polynomials corresponding to assertions outside  $\mathcal{I}(\mathcal{G})$  cannot be identically zero, because they have at least one counterexample. Thus, the set of distributions  $P$ , which exhibit any one of these “spurious” independence assertions, has measure zero. The set of distributions that do not satisfy  $\mathcal{I}(P) = \mathcal{I}(\mathcal{G})$  is the union of these separate sets, one for each spurious independence assertion. The union of a finite number of sets of measure zero is a set of measure zero, proving the result. ■



These results state that for almost all parameterizations  $P$  of the graph  $\mathcal{G}$  (that is, for almost all possible choices of CPDs for the variables), the d-separation test precisely characterizes the independencies that hold for  $P$ . In other words, even if we have a distribution  $P$  that satisfies more independencies than  $\mathcal{I}(\mathcal{G})$ , a slight perturbation of the CPDs of  $P$  will almost always eliminate these “extra” independencies. This guarantee seems to state that such independencies are always accidental, and we will never encounter them in practice. However, as we illustrate in example 3.7, there are cases where our CPDs have certain local structure that is not accidental, and that implies these additional independencies that are not detected by d-separation.

### 3.3.3 An Algorithm for d-Separation

The notion of d-separation allows us to infer independence properties of a distribution  $P$  that factorizes over  $\mathcal{G}$  simply by examining the connectivity of  $\mathcal{G}$ . However, in order to be useful, we need to be able to determine d-separation effectively. Our definition gives us a constructive solution, but a very inefficient one: We can enumerate all trails between  $X$  and  $Y$ , and check each one to see whether it is active. The running time of this algorithm depends on the number of trails in the graph, which can be exponential in the size of the graph.

Fortunately, there is a much more efficient algorithm that requires only linear time in the size of the graph. The algorithm has two phases. We begin by traversing the graph bottom up, from the leaves to the roots, marking all nodes that are in  $Z$  or that have descendants in  $Z$ . Intuitively, these nodes will serve to enable v-structures. In the second phase, we traverse breadth-first from  $X$  to  $Y$ , stopping the traversal along a trail when we get to a blocked node. A node is blocked if: (a) it is the “middle” node in a v-structure and unmarked in phase I, or (b) is not such a node and is in  $Z$ . If our breadth-first search gets us from  $X$  to  $Y$ , then there is an active trail between them.

The precise algorithm is shown in algorithm 3.1. The first phase is straightforward. The second phase is more subtle. For efficiency, and to avoid infinite loops, the algorithm must keep track of all nodes that have been visited, so as to avoid visiting them again. However, in graphs with loops (multiple trails between a pair of nodes), an intermediate node  $Y$  might be involved in several trails, which may require different treatment within the algorithm:

#### Example 3.4

Consider the Bayesian network of figure 3.6, where our task is to find all nodes reachable from  $X$ . Assume that  $Y$  is observed, that is,  $Y \in Z$ . Assume that the algorithm first encounters  $Y$  via the direct edge  $Y \rightarrow X$ . Any extension of this trail is blocked by  $Y$ , and hence the algorithm stops the traversal along this trail. However, the trail  $X \leftarrow Z \rightarrow Y \leftarrow W$  is not blocked by  $Y$ . Thus, when we encounter  $Y$  for the second time via the edge  $Z \rightarrow Y$ , we should not ignore it. Therefore, after

**Algorithm 3.1 Algorithm for finding nodes reachable from  $X$  given  $Z$  via active trails**

```

Procedure Reachable (
   $\mathcal{G}$ , // Bayesian network graph
   $X$ , // Source variable
   $Z$  // Observations
)
1 // Phase I: Insert all ancestors of  $Z$  into  $V$ 
2  $L \leftarrow Z$  // Nodes to be visited
3  $A \leftarrow \emptyset$  // Ancestors of  $Z$ 
4 while  $L \neq \emptyset$ 
5   Select some  $Y$  from  $L$ 
6    $L \leftarrow L - \{Y\}$ 
7   if  $Y \notin A$  then
8      $L \leftarrow L \cup \text{Pa}_Y$  //  $Y$ 's parents need to be visited
9      $A \leftarrow A \cup \{Y\}$  //  $Y$  is ancestor of evidence
10
11 // Phase II: traverse active trails starting from  $X$ 
12  $L \leftarrow \{(X, \uparrow)\}$  // (Node,direction) to be visited
13  $V \leftarrow \emptyset$  // (Node,direction) marked as visited
14  $R \leftarrow \emptyset$  // Nodes reachable via active trail
15 while  $L \neq \emptyset$ 
16   Select some  $(Y, d)$  from  $L$ 
17    $L \leftarrow L - \{(Y, d)\}$ 
18   if  $(Y, d) \notin V$  then
19     if  $Y \notin Z$  then
20        $R \leftarrow R \cup \{Y\}$  //  $Y$  is reachable
21        $V \leftarrow V \cup \{(Y, d)\}$  // Mark  $(Y, d)$  as visited
22       if  $d = \uparrow$  and  $Y \notin Z$  then // Trail up through  $Y$  active if  $Y$  not in  $Z$ 
23         for each  $Z \in \text{Pa}_Y$ 
24            $L \leftarrow L \cup \{(Z, \uparrow)\}$  //  $Y$ 's parents to be visited from bottom
25         for each  $Z \in \text{Ch}_Y$ 
26            $L \leftarrow L \cup \{(Z, \downarrow)\}$  //  $Y$ 's children to be visited from top
27       else if  $d = \downarrow$  then // Trails down through  $Y$ 
28         if  $Y \notin Z$  then
29           // Downward trails to  $Y$ 's children are active
30           for each  $Z \in \text{Ch}_Y$ 
31              $L \leftarrow L \cup \{(Z, \downarrow)\}$  //  $Y$ 's children to be visited from top
32         if  $Y \in A$  then // v-structure trails are active
33           for each  $Z \in \text{Pa}_Y$ 
34              $L \leftarrow L \cup \{(Z, \uparrow)\}$  //  $Y$ 's parents to be visited from bottom
35 return  $R$ 

```

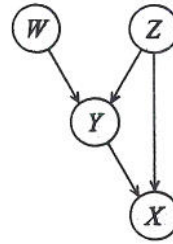


Figure 3.6 A simple example for the d-separation algorithm

the first visit to  $Y$ , we can mark it as visited for the purpose of trails coming in from children of  $Y$ , but not for the purpose of trails coming in from parents of  $Y$ . ■

In general, we see that, for each node  $Y$ , we must keep track separately of whether it has been visited from the top and whether it has been visited from the bottom. Only when both directions have been explored is the node no longer useful for discovering new active trails.

Based on this intuition, we can now show that the algorithm achieves the desired result:

**Theorem 3.6**

The algorithm  $\text{Reachable}(\mathcal{G}, X, Z)$  returns the set of all nodes reachable from  $X$  via trails that are active in  $\mathcal{G}$  given  $Z$ .

The proof is left as an exercise (exercise 3.14).

### 3.3.4 I-Equivalence

The notion of  $\mathcal{I}(\mathcal{G})$  specifies a set of conditional independence assertions that are associated with a graph. This notion allows us to abstract away the details of the graph structure, viewing it purely as a specification of independence properties. In particular, one important implication of this perspective is the observation that very different BN structures can actually be equivalent, in that they encode precisely the same set of conditional independence assertions. Consider, for example, the three networks in figure 3.5a,(b),(c). All three of them encode precisely the same independence assumptions:  $(X \perp Y \mid Z)$ .

**Definition 3.9**

I-equivalence

Two graph structures  $\mathcal{K}_1$  and  $\mathcal{K}_2$  over  $\mathcal{X}$  are I-equivalent if  $\mathcal{I}(\mathcal{K}_1) = \mathcal{I}(\mathcal{K}_2)$ . The set of all graphs over  $\mathcal{X}$  is partitioned into a set of mutually exclusive and exhaustive I-equivalence classes, which are the set of equivalence classes induced by the I-equivalence relation. ■

Note that the v-structure network in figure 3.5d induces a very different set of d-separation assertions, and hence it does not fall into the same I-equivalence class as the first three. Its I-equivalence class contains only that single network.



I-equivalence of two graphs immediately implies that any distribution  $P$  that can be factorized over one of these graphs can be factorized over the other. Furthermore, there is no intrinsic property of  $P$  that would allow us to associate it with one graph rather than an equivalent one. This observation has important implications with respect to our ability to determine the directionality of influence. In particular, although we can determine, for a distribution



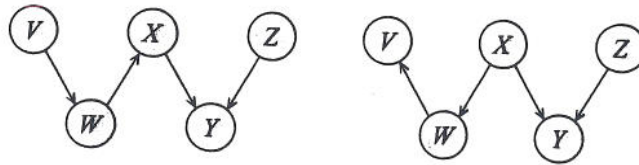


Figure 3.7 Skeletons and v-structures in a network. The two networks shown have the same skeleton and v-structures ( $X \rightarrow Y \leftarrow Z$ ).

$P(X, Y)$ , whether  $X$  and  $Y$  are correlated, there is nothing in the distribution that can help us determine whether the correct structure is  $X \rightarrow Y$  or  $Y \rightarrow X$ . We return to this point when we discuss the causal interpretation of Bayesian networks in chapter 21.

The d-separation criterion allows us to test for I-equivalence using a very simple graph-based algorithm. We start by considering the trails in the networks.

**Definition 3.10**  
skeleton

The skeleton of a Bayesian network graph  $\mathcal{G}$  over  $\mathcal{X}$  is an undirected graph over  $\mathcal{X}$  that contains an edge  $\{X, Y\}$  for every edge  $(X, Y)$  in  $\mathcal{G}$ . ■

In the networks of figure 3.7, the networks (a) and (b) have the same skeleton.

If two networks have a common skeleton, then the set of trails between two variables  $X$  and  $Y$  is same in both networks. If they do not have a common skeleton, we can find a trail in one network that does not exist in the other and use this trail to find a counterexample for the equivalence of the two networks.

Ensuring that the two networks have the same trails is clearly not enough. For example, the networks in figure 3.5 all have the same skeleton. Yet, as the preceding discussion shows, the network of figure 3.5d is not equivalent to the networks of figure 3.5a–(c). The difference, is of course, the v-structure in figure 3.5d. Thus, it seems that if the two networks have the same skeleton and exactly the same set of v-structures, they are equivalent. Indeed, this property provides a sufficient condition for I-equivalence:

**Theorem 3.7**

Let  $\mathcal{G}_1$  and  $\mathcal{G}_2$  be two graphs over  $\mathcal{X}$ . If  $\mathcal{G}_1$  and  $\mathcal{G}_2$  have the same skeleton and the same set of v-structures then they are I-equivalent.

The proof is left as an exercise (see exercise 3.16).

Unfortunately, this characterization is not an equivalence: there are graphs that are I-equivalent but do not have the same set of v-structures. As a counterexample, consider *complete* graphs over a set of variables. Recall that a complete graph is one to which we cannot add additional arcs without causing cycles. Such graphs encode the empty set of conditional independence assertions. Thus, any two complete graphs are I-equivalent. Although they have the same skeleton, they invariably have different v-structures. Thus, by using the criterion on theorem 3.7, we can conclude (in certain cases) only that two networks are I-equivalent, but we cannot use it to guarantee that they are not.

We can provide a stronger condition that does correspond exactly to I-equivalence. Intuitively, the unique independence pattern that we want to associate with a v-structure  $X \rightarrow Z \leftarrow Y$  is that  $X$  and  $Y$  are independent (conditionally on their parents), but dependent given  $Z$ . If there

is a direct edge between  $X$  and  $Y$ , as there was in our example of the complete graph, the first part of this pattern is eliminated.

**Definition 3.11**

immorality

covering edge

A  $v$ -structure  $X \rightarrow Z \leftarrow Y$  is an immorality if there is no direct edge between  $X$  and  $Y$ . If there is such an edge, it is called a covering edge for the  $v$ -structure. ■

Note that not every  $v$ -structure is an immorality, so that two networks with the same immoralities do not necessarily have the same  $v$ -structures. For example, two different complete directed graphs always have the same immoralities (none) but different  $v$ -structures.

**Theorem 3.8**

Let  $\mathcal{G}_1$  and  $\mathcal{G}_2$  be two graphs over  $\mathcal{X}$ . Then  $\mathcal{G}_1$  and  $\mathcal{G}_2$  have the same skeleton and the same set of immoralities if and only if they are I-equivalent.

The proof of this (more difficult) result is also left as an exercise (see exercise 3.17).

We conclude with a final characterization of I-equivalence in terms of local operations on the graph structure.

**Definition 3.12**

covered edge

**Theorem 3.9**

An edge  $X \rightarrow Y$  in a graph  $\mathcal{G}$  is said to be covered if  $\text{Pa}_Y^{\mathcal{G}} = \text{Pa}_X^{\mathcal{G}} \cup \{X\}$ . ■

Two graphs  $\mathcal{G}$  and  $\mathcal{G}'$  are I-equivalent if and only if there exists a sequence of networks  $\mathcal{G} = \mathcal{G}_1, \dots, \mathcal{G}_k = \mathcal{G}'$  that are all I-equivalent to  $\mathcal{G}$  such that the only difference between  $\mathcal{G}_i$  and  $\mathcal{G}_{i+1}$  is a single reversal of a covered edge.

The proof of this theorem is left as an exercise (exercise 3.18).

### 3.4 From Distributions to Graphs

In the previous sections, we showed that, if  $P$  factorizes over  $\mathcal{G}$ , we can derive a rich set of independence assertions that hold for  $P$  by simply examining  $\mathcal{G}$ . This result immediately leads to the idea that we can use a graph as a way of revealing the structure in a distribution. In particular, we can test for independencies in  $P$  by constructing a graph  $\mathcal{G}$  that represents  $P$  and testing d-separation in  $\mathcal{G}$ . As we will see, having a graph that reveals the structure in  $P$  has other important consequences, in terms of reducing the number of parameters required to specify or learn the distribution, and in terms of the complexity of performing inference on the network.

In this section, we examine the following question: Given a distribution  $P$ , to what extent can we construct a graph  $\mathcal{G}$  whose independencies are a reasonable surrogate for the independencies in  $P$ ? It is important to emphasize that we will never actually take a fully specified distribution  $P$  and construct a graph  $\mathcal{G}$  for it: As we discussed, a full joint distribution is much too large to represent explicitly. However, answering this question is an important conceptual exercise, which will help us later on when we try to understand the process of constructing a Bayesian network that represents our model of the world, whether manually or by learning from data.

#### 3.4.1 Minimal I-Maps

One approach to finding a graph that represents a distribution  $P$  is simply to take any graph that is an I-map for  $P$ . The problem with this naive approach is clear: As we saw in example 3.3, the

**Algorithm 3.2 Procedure to build a minimal I-map given an ordering**


---

```

Procedure Build-Minimal-I-Map (
   $X_1, \dots, X_n$  // an ordering of random variables in  $\mathcal{X}$ 
   $\mathcal{I}$  // Set of independencies
)
1  Set  $\mathcal{G}$  to an empty graph over  $\mathcal{X}$ 
2  for  $i = 1, \dots, n$ 
3     $U \leftarrow \{X_1, \dots, X_{i-1}\}$  //  $U$  is the current candidate for parents of  $X_i$ 
4    for  $U' \subseteq \{X_1, \dots, X_{i-1}\}$ 
5      if  $U' \subset U$  and  $(X_i \perp \{X_1, \dots, X_{i-1}\} - U' \mid U') \in \mathcal{I}$  then
6         $U \leftarrow U'$ 
7      // At this stage  $U$  is a minimal set satisfying  $(X_i \perp$ 
8         $\{X_1, \dots, X_{i-1}\} - U \mid U)$ 
9      // Now set  $U$  to be the parents of  $X_i$ 
10     for  $X_j \in U$ 
11       Add  $X_j \rightarrow X_i$  to  $\mathcal{G}$ 
12  return  $\mathcal{G}$ 

```

---

complete graph is an I-map for any distribution, yet it does not reveal any of the independence structure in the distribution. However, examples such as this one are not very interesting. The graph that we used as an I-map is clearly and trivially unrepresentative of the distribution, in that there are edges that are obviously redundant. This intuition leads to the following definition, which we also define more broadly:

**Definition 3.13**  
minimal I-map

A graph  $\mathcal{K}$  is a minimal I-map for a set of independencies  $\mathcal{I}$  if it is an I-map for  $\mathcal{I}$ , and if the removal of even a single edge from  $\mathcal{K}$  renders it not an I-map. ■

This notion of an I-map applies to multiple types of graphs, both Bayesian networks and other types of graphs that we will encounter later on. Moreover, because it refers to a set of independencies  $\mathcal{I}$ , it can be used to define an I-map for a distribution  $P$ , by taking  $\mathcal{I} = \mathcal{I}(P)$ , or to another graph  $\mathcal{K}'$ , by taking  $\mathcal{I} = \mathcal{I}(\mathcal{K}')$ .

Recall that definition 3.5 defines a Bayesian network to be a distribution  $P$  that factorizes over  $\mathcal{G}$ , thereby implying that  $\mathcal{G}$  is an I-map for  $P$ . It is standard to restrict the definition even further, by requiring that  $\mathcal{G}$  be a minimal I-map for  $P$ .

How do we obtain a minimal I-map for the set of independencies induced by a given distribution  $P$ ? The proof of the factorization theorem (theorem 3.1) gives us a procedure, which is shown in algorithm 3.2. We assume we are given a predetermined *variable ordering*, say,  $\{X_1, \dots, X_n\}$ . We now examine each variable  $X_i$ ,  $i = 1, \dots, n$  in turn. For each  $X_i$ , we pick some minimal subset  $U$  of  $\{X_1, \dots, X_{i-1}\}$  to be  $X_i$ 's parents in  $\mathcal{G}$ . More precisely, we require that  $U$  satisfy  $(X_i \perp \{X_1, \dots, X_{i-1}\} - U \mid U)$ , and that no node can be removed from  $U$  without violating this property. We then set  $U$  to be the parents of  $X_i$ .

The proof of theorem 3.1 tells us that, if each node  $X_i$  is independent of  $X_1, \dots, X_{i-1}$  given its parents in  $\mathcal{G}$ , then  $P$  factorizes over  $\mathcal{G}$ . We can then conclude from theorem 3.2 that  $\mathcal{G}$  is an I-map for  $P$ . By construction,  $\mathcal{G}$  is minimal, so that  $\mathcal{G}$  is a minimal I-map for  $P$ .

variable ordering

