

# A Bayesian Method for Learning Belief Networks that Contain Hidden Variables

GREGORY F. COOPER

*Section of Medical Informatics, University of Pittsburgh, B50A Lothrop Hall, Pittsburgh, PA 15261*

**Abstract.** This paper presents a Bayesian method for computing the probability of a Bayesian belief-network structure from a database. In particular, the paper focuses on computing the probability of a belief-network structure that contains a hidden (latent) variable. A hidden variable represents a postulated entity that has not been directly measured. After reviewing related techniques, which previously were reported, this paper presents a new, more efficient method for handling hidden variables in belief networks.

**Keywords:** probabilistic networks, Bayesian belief networks, hidden variables, machine learning, induction

## 1. Introduction

A belief-network structure can provide insight into probabilistic dependencies that exist among the variables in a database. One application is the automated discovery of dependency relationships. The computer program searches for a belief-network structure that has a high posterior probability given the database, and outputs the structure and its probability. A related task is computer-assisted hypothesis testing: The user enters a hypothetical structure of the dependency relationships among a set of variables, and the program calculates the probability of the structure given a database of cases on the variables. We will emphasize this task in the current paper. We also note that given a belief-network structure and a database, we can construct a belief network and use it for computer-based diagnosis and prediction.

In Cooper & Herskovits (1992) a Bayesian method is presented for computing the probability of a belief-network structure that contains hidden variables. The major problem with this approach is that its computational time complexity is exponential in the number of database cases (i.e., records or instances). In this paper, we present a method that has a time complexity that is polynomial in the number of cases, yet it always computes the same result as the exponential method. By being more efficient, this polynomial-time method permits us to solve modeling problems that previously were not computationally feasible. In some instances, which we explain, the polynomial degree of the new method can be high. As the polynomial degree increases, the number of cases considered in the database must decrease in order to maintain computational tractability. We will initially focus our discussion on handling a single hidden variable, and then extend the method to handle multiple hidden variables.

## 2. Background on belief networks

A belief-network *structure*  $B_S$  is a directed acyclic graph in which nodes represent domain variables and arcs between nodes represent probabilistic dependencies (Cooper, 1989; Pearl, 1988). The representation of conditional dependence and independence among variables is the essential function of belief networks. For a detailed discussion of the formal mathematical properties of belief networks, see Pearl (1988).

A belief-network structure,  $B_S$ , is augmented by conditional probabilities,  $B_P$ , to form a belief network  $B$ .<sup>1</sup> Thus,  $B = (B_S, B_P)$ . For each node<sup>2</sup> in a belief-network structure, there is a conditional-probability function that relates this node to its immediate predecessors (parents). We shall use  $\pi_i$  to denote the parent nodes of variable  $x_i$ . If a node has no parents, then a prior-probability function,  $P(x_i)$ , is specified.

Belief networks are capable of representing the probabilities over any discrete sample space: The probability of any sample point in that space can be computed from the probabilities in the belief network. As mentioned, the key feature of belief networks is their explicit representation of the conditional independence and dependence among events. In particular, investigators have shown (Kiiveri, et al., 1984; Pearl, 1988; Shachter, 1986) that the joint probability of any particular instantiation<sup>3</sup> of all  $n$  variables in a belief network can be calculated as follows:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \pi_i),$$

where  $X_i$  represents the instantiation of variable  $x_i$  and  $\pi_i$  represents the instantiation of the parents of  $x_i$ .

Therefore, the joint probability of any instantiation of all the variables in a belief network can be computed as the product of only  $n$  probabilities. In principle, we can recover the complete joint-probability space from the belief-network representation by calculating the joint probabilities that result from every possible instantiation of the  $n$  variables in the network. Thus, we have sufficient information to determine any probability of the form  $P(W | V)$ , where  $W$  and  $V$  are sets of variables with known values (i.e., instantiated variables).

In the last few years, researchers have made significant progress in formalizing the theory of belief networks (Neapolitan, 1990; Pearl, 1988), and in developing more efficient algorithms for probabilistic inference on belief networks (Henrion, 1990). The feasibility of using belief networks in constructing diagnostic systems has been demonstrated in several domains (Agogino & Rege, 1987; Andreassen, et al., 1987; Beinlich, et al., 1989; Chavez & Cooper, 1990; Cooper, 1984; Heckerman, et al., 1992; Henrion & Cooley, 1987; Holtzman, 1989; Suermondt & Amylon, 1989).

Although researchers have made substantial advances in developing the theory and application of belief networks, the actual construction of these networks often remains a difficult, time-consuming task. The task is time-consuming because typically it must be performed manually by an expert or with the help of an expert. In domains that are large or in which there are few, if any, readily available experts, automated methods are needed for augmenting the manual expert-based methods of knowledge acquisition for belief-network

construction. In this paper, we present one such method and show how it can handle hidden variables.

Researchers have developed non-Bayesian methods for discovering the presence of hidden variables in belief-network structures (Spirtes & Glymour, 1990; Spirtes, et al., 1993; Verma & Pearl, 1990). These methods use statistical tests of independence to infer a set of independence and dependence relationships among the model variables. For example, in the case of discrete variable models, the chi-square test has been used (Spirtes, et al., 1993). The algorithms then search to discover a class of belief-network structures that is consistent with the inferred relationships. Thus, the validity of the output of these algorithms depends on whether the conditional independence and dependence relationships among variables can be determined accurately from the available sample of cases. The application of non-Bayesian methods has produced promising results (Spirtes, et al., 1993). Also, in the case of linear systems, non-Bayesian algorithms exist for uncovering relationships between hidden variables and continuous measured variables, as well as uncovering relationships among the hidden variables themselves (Spirtes, et al., 1993).

The non-Bayesian methods do not provide a probability of a belief-network structure that contains one or more hidden variables. Bayesian methods can provide such probabilities. A promising direction for future research may be to use non-Bayesian methods to search for likely alternative belief-network structures that contain hidden variables, and then compare the probabilities of these structures using Bayesian methods. In this paper we discuss an exact Bayesian method for determining the probability of a belief-network structure when there are hidden variables or missing data. We note that in related work researchers have developed simulation (York & Madigan, 1993) and approximation (Spiegelhalter & Cowell, 1992; Spiegelhalter & Lauritzen, 1990) techniques to perform Bayesian parameter estimation in belief-networks when there is missing data.

### 3. An example

In this section we introduce a simple example that we will use throughout the paper to illustrate basic concepts. Assume we have the database shown in Table 1, which we designate as database  $D$ . To be more concrete, suppose  $x_1$  reg's the cause of a disease, and  $x_2$  and  $x_3$  represent two patient symptoms. Assume from independent knowledge that we are quite certain neither symptom causally leads to the other symptom. More particularly, suppose we believe that the causal relationships among the three variables correspond to one of the following two possibilities: (1) the symptom variables are conditionally independent given the state of the disease-etiology variable (see Figure 1), or (2) the symptom variables are caused by an intermediate pathophysiological process between the disease-etiology variable and the two symptoms (see Figure 2). Figures 1 and 2 show the two possible belief-network structures under consideration, which we denote as  $B_{S_1}$  and  $B_{S_2}$ . In  $B_{S_2}$  the variable  $h$  represents an unknown pathophysiological process. In general, for a belief network with this structure,  $x_2$  and  $x_3$  will not be conditionally independent given  $x_1$ . To keep the example simple, we will assume that  $h$  is a binary variable and that we have the following prior probabilities for  $B_{S_1}$  and  $B_{S_2}$ :  $P(B_{S_1}) = P(B_{S_2}) = 0.5$ .

Table 1. A database example, which we denote as  $D$ . The term *case* in the first column denotes a single training instance (record) in the database—as for example, a patient case.

Case	Variable values for each case		
	$x_1$	$x_2$	$x_3$
1	present	absent	present
2	present	present	present
3	absent	absent	absent
4	present	present	present
5	absent	present	absent
6	absent	present	present
7	present	present	present
8	present	present	present
9	absent	present	present
10	present	absent	absent

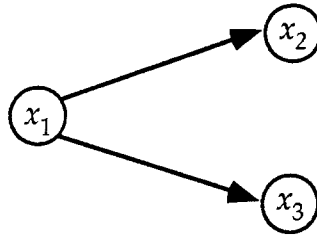


Figure 1. Belief-network structure  $B_{S_1}$ , which represents that  $x_2$  and  $x_3$  are conditionally independent given  $x_1$ .

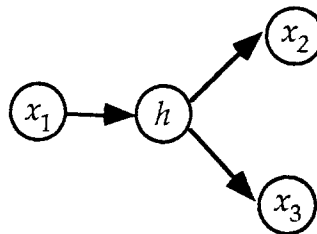


Figure 2. Belief-network structure  $B_{S_2}$ , which represents that  $x_2$  and  $x_3$  are conditionally independent given some hidden variable  $h$ .

Although the current example is hypothetical, Henrion has documented a case in which an expert (who was building a belief network to diagnosis particular apple-tree diseases) first constructed the subnetwork structure shown in Figure 1 (Henrion, 1989). The expert later realized that the belief-network structure in Figure 2 more accurately reflects the relationships among the variables. The methods we propose in this paper are intended to automatically suggest refinements like the one in this example, as well as other refinements.

Our goal is to determine  $P(B_{S_1} | D)$  and  $P(B_{S_2} | D)$ . Using the assumptions and techniques to be presented in Section 4, we can show that  $P(B_{S_1} | D) = 0.25$  and  $P(B_{S_2} | D) = 0.75$ . Deriving these results with the methods in Section 4, however, requires that we perform 1024 operations on the data in Table 1. In Section 5, we introduce a more efficient method that derives the same results in only 240 operations. In Sections 4 and 5 we provide a general analysis of the computational time complexity of both methods.

#### 4. Previous results

Previously, a method was described for Bayesian belief-network-structure learning from a database that is generated by a process that can be accurately modeled by a belief network with no missing values or hidden variables (Cooper & Herskovits, 1992). In Section 4.1, we present those results. In Section 4.2, we show previous extensions of this method that learn belief-network structures containing hidden variables.

##### 4.1. Results when there are no missing data or hidden variables

Suppose for the moment that we have a method for calculating  $P(B_{S_i}, D)$  for some belief-network structure  $B_{S_i}$  and database  $D$ . Let  $Q$  be the set of all those belief-network structures that have a non-zero prior probability. We can derive the posterior probability of  $B_{S_i}$  given  $D$  as  $P(B_{S_i} | D) = P(B_{S_i}, D) / \sum_{B_s \in Q} P(B_{S_i}, D)$ . Sometimes all we want to know is the ratio of the posterior probabilities of two belief-network structures. To calculate such a ratio for belief-network structures  $B_{S_i}$  and  $B_{S_j}$ , we can use the equivalence that  $P(B_{S_i} | D) / P(B_{S_j} | D) = P(B_{S_i}, D) / P(B_{S_j}, D)$ . An algorithm called K2 is reported in Cooper & Herskovits (1992) that makes such comparisons of belief-network structures, as it searches heuristically for the most likely structure.

Our focus in this paper will be on computing the term  $P(B_{S_i}, D)$ , which we can derive using the equivalence that  $P(B_{S_i}, D) = P(D | B_{S_i})P(B_{S_i})$ . The term  $P(B_{S_i})$  represents our prior probability that a process with belief-network structure  $B_{S_i}$  generated data  $D$ . We will assume that  $P(B_{S_i})$  is specified by the user and available. The likelihood term  $P(D | B_{S_i})$  remains to be determined. The following theorem, which is proved in Cooper & Herskovits (1992), provides a method for computing  $P(D | B_S)$ .

**THEOREM 1** *Let  $Z$  be a set of  $n$  discrete variables, where a variable  $x_i$  in  $Z$  has  $r_i$  possible value assignments:  $(v_{i1}, \dots, v_{ir_i})$ . Let  $D$  be a database of  $m$  cases, where each case contains a value assignment for each variable in  $Z$ . Let  $B_S$  denote a belief-network structure containing just the variables in  $Z$ . Each variable  $x_i$  in  $B_S$  has a set of parents,*

which we represent with a list of variables  $\pi_i$ . Let  $w_{ij}$  denote the  $j$ th unique instantiation of  $\pi_i$  relative to  $D$ . Suppose there are  $q_i$  such unique instantiations of  $\pi_i$ . Define  $N_{ijk}$  to be the number of cases in  $D$  in which variable  $x_i$  has the value  $v_{ik}$  and  $\pi_i$  is instantiated as  $w_{ij}$ .

Let  $N_{ij} = \sum_{k=1}^{r_i} N_{ijk}$ . Suppose the following assumptions hold:

1. The variables in  $Z$  are discrete
2. Cases occur independently, given a belief-network model
3. There are no cases that have variables with missing values
4. Before observing  $D$ , we are indifferent regarding which numerical probabilities to assign to the belief network with structure  $B_S$ .

From these four assumptions it follows that

$$P(D | B_S) = \prod_{i=1}^n \prod_{j=1}^{q_i} \frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \prod_{k=1}^{r_i} N_{ijk}! \quad (1)$$

□

The complete proof of Theorem 1 is given in the Appendix in Cooper & Herskovits (1992). We only discuss here the parts of the proof we will need later in this paper. From Assumption 1 we have that

$$P(D | B_S) = \int_{B_P} P(D | B_S, B_P) f(B_P | B_S) dB_P, \quad (2)$$

where  $B_P$  is a vector whose values denote the conditional-probability assignments associated with belief-network structure  $B_S$ , and  $f$  is the conditional-probability-density function over  $B_P$  given  $B_S$ . The integral is over all possible value assignments to  $B_P$ . It follows from the conditional independence of cases expressed in Assumption 2 that Equation 2 can be rewritten as

$$P(D | B_S) = \int_{B_P} \left[ \prod_{t=1}^m P(C_t | B_S, B_P) \right] f(B_P | B_S) dB_P, \quad (3)$$

where  $m$  is the number of cases in  $D$ , and  $C_t$  is the  $t$ th case in  $D$ .

As shown in Cooper & Herskovits (1992), we can apply Assumptions 3 and 4 to derive Equation 1 from Equation 3. Since Equation 3 is all we need for the current paper, we will not show these additional steps here.

As a point of illustration, we describe now one application of Theorem 1. The K2 algorithm mentioned previously uses Equation 1 and assumes uniform priors over belief-network structures. We briefly summarize here the results of applying K2 to data generated by the ALARM belief network. The ALARM network is a laboratory research prototype that models potential anesthesia problems in the operating room (Beinlich, et al. 1989).

ALARM, which contains 37 nodes and 46 arcs, was used to stochastically generate a database of 3000 cases. When given these cases, and a node order, K2 reconstructed the ALARM belief-network structure exactly, except there was one missing arc and one extra arc. The missing arc was not strongly supported by the 3000 cases, and the extra arc was due to the greedy nature of the K2 search algorithm. Additional details about this experiment are discussed in Cooper & Herskovits (1992), and further evaluation of K2 on other databases is presented in Herskovits (1991).

In this paper, we will adopt Assumptions 1, 2, and 4 listed in Theorem 1, although Assumption 4 can be relaxed, as shown in Cooper & Herskovits (1992). We note that Assumption 4 implies a uniform prior over each of the conditional probability distributions (the parameters of the model) associated with a given  $B_S$  (the model structure).

In this paper, we will relax Assumption 3 in one important way: All variables in the belief network are assumed to have known values, except one variable for which we know none of its values. This variable is a hidden variable.

#### 4.2. Previous hidden-variable results

In this section, we consider belief-network structures that contain a single hidden variable, which we denote as  $h$ . In Cooper & Herskovits (1992) we describe a more general (yet still inefficient) method that handles both missing data and multiple hidden variables.

Let  $h$  denote a hidden variable in  $B_S$ . We use  $h_t$  as a variable that represents  $h$  in case  $t$ —we call  $h_t$  a *case-specific hidden variable*. Let  $C_t$  denote the set of variable assignments for those variables in the  $t$ th case that have known values (i.e., all the variables except  $h$ ). The probability of the  $t$ th case can be computed as

$$P(C_t | B_S, B_P) = \sum_{h_t} P(C_t, h_t | B_S, B_P), \quad (4)$$

where the sum  $\sum_{h_t}$  indicates that  $h_t$  is run through all the possible values of  $h$ .

By substituting Equation 4 into Equation 3, we obtain

$$P(D | B_S) = \int_{B_P} \left[ \prod_{t=1}^m \left[ \sum_{h_t} P(C_t, h_t | B_S, B_P) \right] \right] f(B_P | B_S) dB_P. \quad (5)$$

We now can rearrange Equation 5 as follows by converting an integral of a product of a sum into sums of an integral of a product:

$$P(D | B_S) = \sum_{h_1} \cdots \sum_{h_m} \int_{B_P} \left[ \prod_{t=1}^m P(C_t, h_t | B_S, B_P) \right] f(B_P | B_S) dB_P. \quad (6)$$

The sums in Equation 6 are taken over every possible value of the hidden variable for every possible case in the database. Equation 6 is a sum over the type integral that occurs in Equation 3, and as Theorem 1 shows, we can solve such integrals with Equation 1. Thus, Equation 6 can be solved by multiple applications of Equation 1. The complexity

of computing Equation 6 is exponential in the number of cases  $m$  in the database. In particular, if the hidden variable can take on one of  $r$  possible values, then to solve Equation 6 directly requires  $r^m$  applications of Equation 1. This high level of complexity is not computationally tractable for most applications. In Section 5, we describe a more efficient method for computing Equation 6.

*Example:* Let us augment Table 1 to represent a hidden variable. In particular, this new table is the same as Table 1, except that a column is added for the hidden variable  $h$ . For each of the ten cases in the table, the variable  $h$  is represented by a *case-specific hidden variable*  $h_t$ , for  $t = 1$  to 10. Each case-specific hidden variable can be instantiated to either the value *absent* or *present*. We apply Equation 6 to this table in order to compute  $P(D \mid B_{S_2})$ , where  $B_{S_2}$  is the belief-network structure in Figure 2. Since there are ten cases in  $D$ , the total number of instantiations of the ten binary variables  $h_1, h_2, \dots, h_{10}$  is  $2^{10} = 1024$ . Thus, by using Equation 6 to solve for  $P(D \mid B_{S_2})$ , Equation 1 must be solved 1024 times. If we use Equation 6, then as the number of cases with a hidden variable increases, the number of times we must solve Equation 1 grows exponentially. To address this problem of exponential growth in complexity, we introduce a more efficient solution method in the next section.

## 5. A method for handling hidden variables more efficiently

In this section we present a method that computes Equation 6, but does not use exhaustive summation over the values of the case-specific hidden variables. Instead, we aggregate cases in the database into groups, and we use these groups to calculate a result equivalent to one produced by Equation 6. Before introducing the general method, we first illustrate the basic ideas underlying it by showing how to apply the method to the example previously introduced.

### 5.1. Application to the example

In this section we partially formalize the new hidden-variable method by using the previous example to illustrate basic concepts. In Section 5.2 we complete the formalization to the general case.

According to Assumption 2 in Theorem 1, cases are assumed to occur independently given a belief-network model, which implies that the order of the cases in the database is irrelevant. Therefore, in Table 2 we have aggregated identical cases into groups. There are six groups: four of the groups contain only one case, one group contains two cases, and another group contains four cases.

The basic idea underlying a more efficient method for handling hidden variable  $h$  is this: we compute the integrals in Equation 6 using groups of cases rather than individual cases. Recall that the integral in Equation 6 is solved by Equation 1. According to Equation 1, it does not matter *which* case-specific hidden variable in a given group has a particular value. Rather, all that matters is the total number of values of *present* and of *absent* for the



case-specific hidden variables in a group. Thus, for example, the solution of the integral in Equation 6 will be the same, regardless of whether we instantiate the case-specific hidden variables in group 2 as  $(h_6 := present, h_9 := absent)$  or as  $(h_6 := absent, h_9 := present)$ .

Suppose a given group  $G_i$  contains  $c_i$  case-specific hidden variables, and it has  $d_i$  of those variables instantiated to the value *present* and  $c_i - d_i$  of the variables instantiated to *absent*. In this situation, we will say that  $G_i$  is *set* to the value  $d_i$ . Note that  $d_i$  can range from 0 to  $c_i$ , and thus,  $G_i$  has  $c_i + 1$  possible settings. Rather than compute the integral in Equation 6 over all possible instantiations of the case-specific hidden variables in  $D$ , we compute only over the possible settings for each group. Note that there are  $\binom{c_i}{d_i}$  ways to instantiate the case-specific hidden variables in  $G_i$  so that  $d_i$  of them have the value *present* and  $c_i - d_i$  of them have the value *absent*. Suppose there are  $u$  total groups. According to Equation 6, for a given setting of each group, the integral in that equation is solved  $\prod_{i=1,u} \binom{c_i}{d_i}$  times. Thus, for a given setting of the groups, we must multiply our solution of the integral by the factor  $\prod_{i=1,u} \binom{c_i}{d_i}$ .

Consider, as an example, computing the integral in Equation 6 for the situation in which two of the case-specific hidden variables in group 1 have the value *present* (and thus, two must have the value *absent*), one case-specific hidden variable in group 2 has the value *present*, and in groups 3 through 6, each case-specific hidden variable has the value *present*. These settings of the groups correspond to  $\binom{4}{2} \binom{2}{1} \binom{1}{1} \binom{1}{1} \binom{1}{1} \binom{1}{1} = 12$  solutions to the integral in Equation 6, if the exhaustive technique represented by that equation is used. However, each of these 12 solutions will be the same. The solutions are the same because for a given setting of the groups, Equation 1 is insensitive to any particular instantiation of the case-specific hidden variables that is consistent with those given group settings. Thus, to be more efficient (for the given group settings in this example), we solve the integral in Equation 6 only once (by using Equation 1), and we multiply that solution to 12. We perform this procedure for each possible setting of the groups. We sum the results of each such procedure application to obtain a solution equivalent to that obtained by using Equation 6. The total number of joint settings of all the groups is equal the product of the possible settings of each group. For the example, this product is  $5 \times 3 \times 2 \times 2 \times 2 \times 2 = 240$ . Thus, we must solve Equation 1 only 240 times, rather than the 1024 times required by Equation 6.

## 5.2. The general method

In this section we generalize the approach discussed in Section 5.1. We extend hidden variables to be  $r$ -valued discrete variables, for any  $r \geq 2$ . We also provide a general formula for computing  $P(D | B_S)$  that is more efficient than Equation 6. In the next paragraph we establish some terms and definitions that will be used in a theorem that follows.

We will designate the possible values of  $h$  as  $(v_{h1}, \dots, v_{hr_h})$ , where  $r_h \geq 2$ . Let  $MB(h)$  be the Markov blanket (Pearl, 1988) of variable  $h$ .<sup>4</sup> For the variables in  $MB(h)$ , consider some instantiation of them that exists in database  $D$ , and call it instantiation  $i$ . Suppose there are  $u$  unique such instantiations. Consider a set containing all the cases in  $D$  with instantiation  $i$ . Suppose there are  $c_i$  such cases. Let  $G_i$  represent the set of  $c_i$  case-specific hidden variables in these cases. A *setting* of  $G_i$  is an instantiation of the variables in  $G_i$  that is consistent with the distribution  $(y_{i1}, \dots, y_{ir_h})$ , where  $y_{ij}$  designates

Table 2. In this table cases are aggregated into groups that are demarcated by bold horizontal lines. Note in the second column that cases are classified according to the group they are in.

Case	Group	Variable values for each case			
		$x_1$	$h$	$x_2$	$x_3$
2	1	present	$h_2$	present	present
4	1	present	$h_4$	present	present
7	1	present	$h_7$	present	present
8	1	present	$h_8$	present	present
6	2	absent	$h_6$	present	present
9	2	absent	$h_9$	present	present
3	3	absent	$h_3$	absent	absent
10	4	present	$h_{10}$	absent	absent
1	5	present	$h_1$	absent	present
5	6	absent	$h_5$	present	absent

the number of variables in  $G_i$  that are instantiated to the value  $v_{h_j}$ . If an instantiation  $I_a$  of the variables in  $G_i$  has the same distribution  $(y_{i1}, \dots, y_{ir_h})$  as an instantiation  $I_b$ , then  $I_a$  and  $I_b$  correspond to the same setting of  $G_i$ . Let  $\sum_{G_i}$  denote the sum over all the possible settings of  $G_i$ . For each setting of  $G_i$  in this summation, the case-specific hidden variables in  $G_i$  are instantiated to some set of values that correspond to this setting of  $G_i$ . Let  $g(G_1, \dots, G_u)$  denote the total number of possible instantiations of the  $m$  case-specific hidden variables in  $G_1 \cup \dots \cup G_u$  that are consistent with a given joint setting of the groups  $G_1, \dots, G_u$ . Since each  $G_i$  contains a unique set of case-specific variables, it follows that  $g(G_1, \dots, G_u) = g(G_1) \times \dots \times g(G_u)$ . A term  $g(G_i)$  denotes the total number of possible settings of the  $c_i$  variables in  $G_i$ , and as shown in Tucker (1984, page 187, Theorem 1), this term can be computed as

$$g(G_i) = \frac{c_i!}{y_{i1}!y_{i2}! \dots y_{ir_h}!}.$$

**THEOREM 2** *Let  $D$  be a database of  $m$  cases, where each case contains a value assignment for each variable in  $Z$ , except for a hidden variable  $h$ , which has no value assignments. Let  $B_S$  denote a belief-network structure containing just the variables in  $Z$ . Suppose that among the  $m$  cases in  $D$  there are  $u$  unique instantiations of the variables in  $MB(h)$ . Given these conditions and the Assumptions 1, 2 and 4 in Theorem 1, it follows that*

$$P(D | B_S) = \sum_{G_1} \dots \sum_{G_u} g(G_1, \dots, G_u) \int_{B_P} \left[ \prod_{t=1}^m P(C_t, h_t | B_S, B_P) \right] f(B_P | B_S) dB_P. \quad (7)$$

**Proof:** Recall that Equation 1 solves the integral in Equation 6. Note that the result determined by Equation 1 is only sensitive to the values of  $r_i$ ,  $N_{ijk}$  and  $N_{ij}$ . The values of  $N_{ijk}$  and  $N_{ij}$  represent numerical summaries of the database (i.e., counts), which are insensitive to the particular order of cases that produce those counts. Thus, in Equation 6, if one instantiation of  $h_1, \dots, h_m$  leads to the same values of  $N_{ijk}$  and  $N_{ij}$  as another instantiation, then the result from Equation 1 will be the same for both instantiations. A given setting of  $G_1, \dots, G_u$  in Equation 7 corresponds to those instantiations of  $h_1, \dots, h_m$  that produce the same values of  $N_{ijk}$  and  $N_{ij}$  in Equation 1. In particular, for a given setting of  $G_1, \dots, G_u$  there are  $g(G_1, \dots, G_u)$  instantiations of  $h_1, \dots, h_m$  in Equation 6 that produce the same result when the integral in that equation is solved using Equation 1. Thus, to compute  $P(D | B_S)$  using Equation 6, we can solve the integral in that equation using Equation 1 for each setting of  $G_1, \dots, G_u$  and then multiply by  $g(G_1, \dots, G_u)$  the solution of the integral for each such setting. Equation 7 computes  $P(D | B_S)$  in the manner just outlined.  $\square$

### 5.3. Computational time complexity

In this section, we derive the time complexity for computing Equation 7. We will assume that each variable in the database can be assigned any one of at most  $r$  possible discrete values. As shown in Cooper & Herskovits (1992), the complexity of computing the integral in Equation 7 is  $O(mn^2r)$ . Now consider the number of times we must solve the integral in Equation 7. We will use the following lemma.

LEMMA 1 *Let  $x_1, \dots, x_u$  denote real variables, and let  $m$  and  $k$  designate non-negative real constants. Subject to the following constraints: (1)  $x_1 + \dots + x_u = m$ , and (2)  $0 \leq x_1 \leq m, \dots, 0 \leq x_u \leq m$ , the following equivalence is valid:*

$$\max_{x_1, \dots, x_u} \left[ \prod_{i=1}^u (x_i + k) \right] = \prod_{i=1}^u \left( \frac{m}{u} + k \right).$$

*If the variables  $x_1, \dots, x_u$  are further constrained to be integer valued, then the maximization on the left is less than or equal to the product on the right.*  $\square$

THEOREM 3 *The time required to solve Equation 7 is*

$$O((mn^2r)(m/u + r - 1)^{u(r-1)}). \quad (8)$$

**Proof:** The number of settings of  $G_i$  is equivalent to the number of ways of distributing  $c_i$  identical objects into  $r_h$  different boxes, which as shown in Tucker (1984, page 194) is equal to  $(c_i + r_h - 1)/c_i!(r_h - 1)!$ . Thus, the integral in Equation 7 is solved  $\prod_{i=1, u} (c_i + r_h - 1)/c_i!(r_h - 1)!$  times. We can express  $(c_i + r_h - 1)/c_i!(r_h - 1)!$  equivalently as  $(c_i + r_h - 1)(c_i + r_h - 2) \cdots (c_i + 1)/(r_h - 1)!$ , which is no greater than  $(c_i + r_h - 1)^{r_h - 1}$ . Therefore, the integral in Equation 7 is solved no more than  $\prod_{i=1, u} (c_i + r_h - 1)^{r_h - 1} = (\prod_{i=1, u} (c_i + r_h - 1))^{r_h - 1}$  times. By Lemma 1, this last product can be no greater than

$\prod_{i=1,u}(m/u + r_h - 1)$ , and thus the integral in Equation 7 is solved no more than the following number of times:  $(\prod_{i=1,u}(m/u + r_h - 1))^{r_h - 1} = \prod_{i=1,u}(m/u + r_h - 1)^{r_h - 1} = (m/u + r_h - 1)^{u(r_h - 1)}$ . Thus, the number of times the integral in Equation 7 is solved is  $O((m/u + r_h - 1)^{u(r_h - 1)})$ .

Since the integral in Equation 7 is solved  $O((m/u + r_h - 1)^{u(r_h - 1)})$  times and each solution requires,  $O(mn^2r)$  time, a bound on the total time required to solve Equation 7 is given by Equation 8.

Equation 8 does not take into account the time required to create the groups  $G_1, \dots, G_u$ . We can generate these groups by inserting the  $m$  cases into a rooted tree with a branching factor of at most  $r$ , where each path in the tree represents a unique instantiation of the variables in  $MB(h)$ . Such a tree can be constructed in  $O(mr|MB(h)|) = O(mrn)$  time, and thus, Equation 8 still represents the time complexity required to solve Equation 7.  $\square$

Recall that  $u$  denotes only the number of unique instantiations *actually realized* in database  $D$  of the variables in the Markov blanket of hidden variable  $h$ . The number of such unique instantiations significantly influences the efficiency with which we can compute Equation 7. For any finite belief-network, the number of such unique instantiations reaches a maximum, regardless of how many cases there are in the database. Recall that  $r$  denotes the maximum number of possible values for any variable in the database. If  $u$  and  $r$  are bounded from above, then the time to solve Equation 7 is bounded from above by a function that is polynomial in the number of variables  $n$  and the number of cases  $m$ . If  $u$  or  $r$  is large, however, the polynomial will be of high degree.

## 6. Results on an example using simulated data

This section provides an example application of the methods discussed in this paper. We emphasize that this section contains only a single example, rather than a thorough empirical evaluation of the methods in the paper.

We generated a database of 45 cases from a belief network by using a stochastic sampling method (Henrion, 1988). Let us denote this belief network as  $B^*$  and the database as  $D^*$ . The belief-network structure of  $B^*$  is shown in Figure 2. The probabilities of  $B^*$  are as follows:

$$\begin{aligned} P(x_1 = \text{present}) &= 0.6 \\ P(h = \text{present} \mid x_1 = \text{present}) &= 0.9 \\ P(h = \text{present} \mid x_1 = \text{absent}) &= 0.3 \\ P(x_2 = \text{present} \mid h = \text{present}) &= 0.9 \\ P(x_2 = \text{present} \mid h = \text{absent}) &= 0.05 \\ P(x_3 = \text{present} \mid h = \text{present}) &= 0.8 \\ P(x_3 = \text{present} \mid h = \text{absent}) &= 0.05 \end{aligned}$$

Using the 45 cases in  $D^*$ , we computed the probability of the structures in Figures 1 and 2 by applying both the exhaustive method in Section 4.2 (call it EXH) corresponding to

*Table 3.* The time required by the EXH and MB methods to compute the posterior probability of belief-network structure  $B_{S_2}$  shown in Figure 2, using the first 5, 10, 15, . . . , or 45 cases in database  $D^*$ . The symbol s denotes seconds, m denotes minutes, h denotes hours, and y denotes years. Only the first 5, 10, and 15 cases were computed using EXH; the remaining computation times for EXH were made by using a conservative extrapolation that bounds the results from below.

cases:	5	10	15	20	25	30	35	40	45
EXH:	0.1 s	2.4 s	101 s	> 53 m	> 28 h	> 38 d	> 3 y	> 107 y	> 3424 y
MB:	0.1 s	0.2 s	1 s	9 s	21 s	77 s	139 s	232 s	867 s

Equation 6, and the Markov-blanket method in Section 5.2 (call it MB) corresponding to Equation 7. Table 3 shows the timing results of applying EXH and MB to compute the posterior probability of the structure in Figure 2 for the first 5, 10, 15, . . . , and 45 cases in  $D^*$  when using a Macintosh IIfx (with a Daystar cache card) and Think Pascal 4.0. The EXH method could feasibly handle a database with up to about 15 cases on this machine, whereas the MB method could feasibly handle about 40 cases.

These results, while limited in scope, do illustrate two general points. First, the MB method has a run time that in general grows much more slowly than does the run time of EXH. Second, the run time of MB can still grow rapidly, even though it is a much slower growth than that of EXH. In this example, the run time of MB grows polynomially in the number of cases, where the degree of the polynomial is 8, because 8 is the number of unique instantiations in  $D^*$  of the 3 binary variables in the Markov blanket of  $h$  (namely,  $x_1$ ,  $x_2$ , and  $x_3$ ).

Let us assume that belief-network structure  $B_{S_1}$  in Figure 1 and structure  $B_{S_2}$  in Figure 2 each have a prior probability of 0.5. Given these priors, the probability  $P(B_{S_2} | D)$  is plotted in Figure 3 as a function of the number of cases in  $D$ , when  $D$  is taken to be the first 5, 10, 15, . . . , and 45 cases in  $D^*$ . The first 5 cases did not provide enough data to indicate that the most likely data-generating belief network had structure  $B_{S_2}$ . Beginning with the first 10 cases, however, the posterior probability of  $B_{S_2}$  remains greater than that of  $B_{S_1}$ . The dip between case size 10 and 25 in Figure 3 is likely the result of variation due to statistical sampling, which typically is more prominent when the number of samples is not large.

## 7. Extensions and open problems

In this section we briefly discuss several possible extensions to the MB method.

### 7.1. Multiple hidden variables

We can improve the efficiency of computing Equation 7 by moving some terms to the outside of the sums. This becomes more apparent if Equation 7 is rewritten with the integral replaced by its solution (as given by Equation 1):

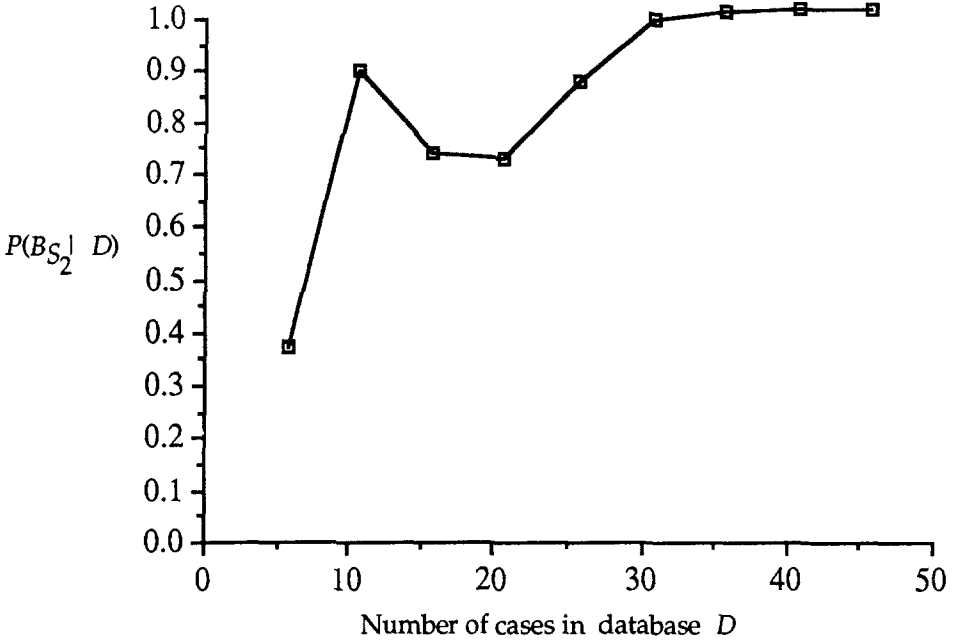


Figure 3. A plot of the posterior probability of belief-network structure  $B_{S_2}$  as  $D$  is set to be the first 5, 10, 15, . . . , and 45 cases in database  $D^*$ . Lines are drawn between data points to highlight changes between those points.

$$P(D | B_S) = \sum_{G_1} \cdots \sum_{G_u} g(G_1, \dots, G_u) \prod_{i=1}^n \prod_{j=1}^{q_i} \frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \prod_{k=1}^{r_i} N_{ijk}!. \quad (9)$$

If a node  $x_i$  is not equal to  $h$  and is not a child of  $h$ , then in Equation 9 the result of the products is not influenced by the settings made in the sums, and thus we can move these products in front of the sums, and calculate them only once. Under these circumstances, Equation 9 becomes

$$P(D | B_S) = \left[ \prod_{i \in Z - C^*(h)} \prod_{j=1}^{q_i} \frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \prod_{k=1}^{r_i} N_{ijk}! \right] \times \left[ \sum_{G_1} \cdots \sum_{G_u} g(G_1, \dots, G_u) \prod_{i \in C^*(h)} \prod_{j=1}^{q_i} \frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \prod_{k=1}^{r_i} N_{ijk}! \right], \quad (10)$$

where  $C^*(h)$  is a set containing just  $h$  and the children of  $h$ . If there is a hidden variable in the set  $Z - C^*(h)$  that is not in the Markov blanket of  $h$ , then we can apply Equation 7 to compute the products that appear in the first set of square brackets in Equation 10. This process of successively moving products out of sums can be repeated for any number

of hidden variables, as long as none of the hidden variables have intersecting Markov boundaries. If the Markov boundaries of two or more hidden variables do have a non-empty intersection, then we can in effect treat these variables like a single hidden variable (of higher dimension), and apply Equation 7 as before. This technique therefore provides a general approach for handling multiple hidden variables. The technique can, however, become computationally demanding, and thus, it is an open problem to discover more efficient methods for handling multiple variables that are hidden or that have some missing data.

### 7.2. *Bounding the space of hidden-variable models*

One difficulty in considering the possibility of hidden variables is that there are an unlimited number of them and thus an unlimited number of belief-network structures that can contain them. This is a difficult, open problem for which we outline here several approaches. One way to address the problem is simply to limit the number of hidden variables in the belief networks that we postulate, and then use a greedy search method similar to K2 (Cooper & Herskovits, 1992) to try to locate a highly probable belief-network structure that contains hidden variables. Another approach is to specify explicitly nonzero priors for only a limited number of belief-network structures that contain hidden variables. This may be possible if we have knowledge that tightly constrains a priori the possible set of structures. One other possible approach is to use non-Bayesian methods to suggest likely locations for hidden variables, as discussed in Spirtes & Glymour (1990), Spirtes, et al. (1993), and Verma & Pearl (1990). We then could begin our search using Bayesian methods by initially postulating hidden variables where the non-Bayesian methods suggest they might exist.

Another problem is to determine the number of values to define for a hidden variable. One approach is to try different numbers of values. That is, we make the number of values of each hidden variable be a parameter in the search space of belief-network structures. We note that some types of unsupervised learning have close parallels to discovering the number of values to assign to hidden variables. For example, researchers have successfully applied unsupervised Bayesian learning methods to determine the most probable number of values of a single, hidden classification variable (Cheeseman, et al., 1988). We believe that similar methods may prove useful in addressing the problem of learning the number of values of hidden variables in belief networks.

### 7.3. *Handling missing values*

Missing values in the database can be viewed as partially hidden variables. We can use a slight variation of Equation 7 to handle missing values for a given variable  $h$  that is partially hidden. Since the number of missing values of a variable is typically much less than the total number of cases, handling missing values will usually be less computationally demanding than handling hidden variables, for the methods discussed in this paper. By Assumption 4 in Theorem 1, for a given belief-network structure  $B_S$ , if a variable  $x_i$  has a missing value in a case, then  $x_i$  is believed to be just as likely to have one value as any other. If this assumption

is not valid, then Assumption 4 should be modified to allow the expression of more general distributions of missing values, as for example, by using Dirichlet distributions (Cooper & Herskovits, 1992).

## Acknowledgments

Gregory Piatetsky-Shapiro and the anonymous reviewers provided helpful comments on an earlier version of this paper.

Support was provided by the National Science Foundation under grant IRI-9111590.

## Appendix

This appendix contains a glossary of terms that appear in the paper.

$D$	A database of cases (i.e., training instances or records).
$m$	The number of cases in $D$ .
$n$	The number of variables in database $D$ .
$B$	A belief network.
$B_S$	A belief-network structure.
$B_P$	A vector whose elements denote the conditional probability assignments associated with a given belief-network structure.
$f(B_P   B_S)$	A conditional-probability density function over $B_P$ given $B_S$ .
$Z$	A set of variables corresponding to the nodes in belief-network $B$ .
$x_i$	A variable in $D$ and a corresponding node in belief network $B$ .
$\pi_i$	The set of parents of $x_i$ in $B$ .
$r_i$	The number of possible value assignments for $x_i$ .
$w_{ij}$	The $j$ th unique instantiation of $\pi_i$ relative to the cases in $D$ .
$q_i$	The number of unique instantiations of $\pi_i$ relative to the cases in $D$ .
$N_{ijk}$	The number of cases in $D$ in which $x_i$ has the value $v_{ik}$ and $\pi_i$ is instantiated as $w_{ij}$ .
$N_{ij}$	A term equal to $\sum_{k=1}^{r_i} N_{ijk}$ .
$h$	A hidden variable. That is, a variable for which there are no measured values in $D$ .
$h_t$	A hidden variable specific to case $t$ .
$G_i$	A set of case-specific hidden variables. See Section 5.2.
$g(G_1, \dots, G_u)$	The total number of possible instantiations of the $m$ case-specific hidden variables in $G_1 \cup \dots \cup G_u$ that are consistent with the joint setting of the groups $G_1, \dots, G_u$ .

## Notes

1. The Appendix contains a glossary of terms used in this paper.
2. Since there is a one-to-one correspondence between a node in  $B_S$  and variable in  $B_P$ , we shall use the terms *node* and *variable* interchangeably.



3. An instantiated variable is a variable with an assigned value.
4. A Markov blanket of a node  $x_i$  is given by the set of nodes consisting of the parents of  $x_i$ , the children of  $x_i$ , and all the parents of the children of  $x_i$  except  $x_i$  itself. Knowing the values of the nodes in the Markov blanket of  $x_i$  makes the probability distribution over  $x_i$  independent of the values of all nodes outside the Markov blanket of  $x_i$ .

## References

- Agogino, A.M. and Rege, A. (1987). "IDES: Influence diagram based expert system," *Mathematical Modelling* Vol. 8, pp. 227–233.
- Andreassen, S., Woldbye, M., Falck, B. and Andersen, S.K. (1987). "MUNIN—A causal probabilistic network for interpretation of electromyographic findings," in *Proceedings of the International Joint Conference on Artificial Intelligence*, Milan, Italy, pp. 366–372.
- Beinlich, I.A., Suermondt, H.J., Chavez, R.M. and Cooper, G.F. (1989). "The ALARM monitoring system: A case study with two probabilistic inference techniques for belief networks," in *Proceedings of the Second European Conference on Artificial Intelligence in Medicine*, London, England, pp. 247–256.
- Chavez, R.M. and Cooper, G.F. (1990). KNET: Integrating hypermedia and normative Bayesian modeling. In: Shachter, R.D., Levitt, T.S., Kanal, L.N. and Lemmer, J.F. (Eds.), *Uncertainty in Artificial Intelligence 4*. North-Holland, Amsterdam, pp. 339–349.
- Cheeseman, P., Self, M., Kelly, J., Taylor, W., Freeman, D. and Stutz, J. (1988). "Bayesian classification," in *Proceedings of AAAI*, St. Paul, MN, pp. 607–611.
- Cooper, G.F. (1984). *NESTOR: A Computer-Based Medical Diagnostic Aid that Integrates Causal and Probabilistic Knowledge*, Doctoral dissertation, Medical Information Sciences, Stanford University.
- Cooper, G.F. (1989). "Current research directions in the development of expert systems based on belief networks," *Applied Stochastic Models and Data Analysis* Vol. 5, pp. 39–52.
- Cooper, G.F. and Herskovits, E. (1992). "A Bayesian method for the induction of probabilistic networks from data," *Machine Learning* Vol. 9, pp. 309–347.
- Heckerman, D.E., Horvitz, E.J. and Nathwani, B.N. (1992). "Toward normative expert systems: Part I. The Pathfinder project," Vol. 31, pp. 90–105.
- Henrion, M. (1988). Propagating uncertainty in Bayesian networks by logic sampling. In: Lemmer, J.F. and Kanal, L.N. (Eds.), *Uncertainty in Artificial Intelligence 2*. North-Holland, Amsterdam, pp. 149–163.
- Henrion, M. (1989). Some practical issues in constructing belief networks. In: Kanal, L.N., Levitt, T.S. and Lemmer, J.F. (Eds.), *Uncertainty in Artificial Intelligence 3*. Elsevier, New York, pp. 161–173.
- Henrion, M. (1990). An introduction to algorithms for inference in belief nets. In: Henrion, M., Shachter, R.D., Kanal, L.N. and Lemmer, J.F. (Eds.), *Uncertainty in Artificial Intelligence 5*. North-Holland, Amsterdam, pp. 129–138.
- Henrion, M. and Cooley, D.R. (1987). "An experimental comparison of knowledge engineering for expert systems and for decision analysis," in *Proceedings of AAAI*, Seattle, WA, pp. 471–476.
- Herskovits, E.H. (1991). *Computer-Based Probabilistic-Network Construction*, Doctoral dissertation, Medical Information Sciences, Stanford University.
- Holtzman, S. (1989). *Intelligent Decision Systems*, Addison-Wesley, Reading, MA.
- Kiiveri, H., Speed, T.P. and Carlin, J.B. (1984). "Recursive causal models," *Journal of the Australian Mathematical Society* Vol. 36, pp. 30–52.
- Neapolitan, R. (1990). *Probabilistic Reasoning in Expert Systems*, John Wiley & Sons, New York.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems*, Morgan Kaufmann, San Mateo, CA.
- Shachter, R.D. (1986). Intelligent probabilistic inference. In: Kanal, L.N. and Lemmer, J.F. (Eds.), *Uncertainty in Artificial Intelligence*. North-Holland, Amsterdam, pp. 371–382.
- Spiegelhalter, D.J. and Cowell, R.G. (1992). Learning in probabilistic expert systems. In: Bernardo, J.M., Berger, J.O., Dawid, A.P. and Smith, A.F.M. (Eds.), *Bayesian Statistics 4*. Oxford University Press, Oxford, pp. 1–17.
- Spiegelhalter, D.J. and Lauritzen, S.L. (1990). "Sequential updating of conditional probabilities on directed graphical structures," *Networks* Vol. 20, pp. 579–606.
- Spirtes, P. and Glymour, C. (1990). "Causal structure among measured variables preserved with unmeasured variables," Report CMU-LCL-90-5, Department of Philosophy, Carnegie-Mellon University.
- Spirtes, P., Glymour, C. and Scheines, R. (1993). *Causation, Prediction, and Search*, Springer-Verlag, New York.

- Suermondt, H.J. and Amylon, M.D. (1989). "Probabilistic prediction of the outcome of bone-marrow transplantation," in *Proceedings of the Symposium on Computer Applications in Medical Care*, Washington, DC, pp. 208–212.
- Tucker, A. (1984). *Applied Combinatorics*, John Wiley & Sons, New York.
- Verma, T.S. and Pearl, J. (1990). "Equivalence and synthesis of causal models," in *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, Cambridge, Massachusetts, pp. 220–227.
- York, J. and Madigan, D. (1993). "Markov chain Monte Carlo methods for hierarchical Bayesian expert systems," in *Proceedings of the Conference on Artificial Intelligence in Statistics*, Orlando, Florida, pp. 433–439.