

Lecture 8: April 14, 2005

*Lecturer: Irit Gat-Viks**Scribe: Uri Avni and Liza Potikha¹*

8.1 Exploiting Structure in Probability Distributions

8.1.1 Introduction

A major goal in biology is to model biological processes given biological data. It is essential for such a model to be of probabilistic nature, since measurements are noisy and molecular biology itself has stochastic characteristics. The problem is that biological processes are composed of a large number of parameters and representing their probabilistic properties requires an exponential, thus very large, number of variables. In this lecture Bayesian networks [7, 19] will be introduced². These networks take advantage of the local dependence occurring in biological processes in order to represent their probability distribution in a compact way. These networks are useful for describing systems composed of locally interacting components, that is, the value of each component depends on the values of a small number of other components. Bayesian networks also provide us a graphical, easy to understand, model of causal influence, as will be discussed later.

One advantage in using Bayesian networks is the fact that the statistical foundations for learning Bayesian networks from observations, and computational algorithms to do so are well understood and have been used successfully for many applications. These applications include medical and fault diagnosis expert systems, monitoring systems, information access technologies, speech recognition systems, and finally analysis and classification for biological sequencing.

This lecture starts with a probability primer, describing the important definitions and results needed to understand Bayesian networks. Later, Bayesian networks are introduced and an example to their compactness of representation is given. Methods to compute probability distributions based on a Bayesian network are then presented. Once the Bayesian networks were described as a compact way to represent a stochastic (local) system, methods to deduce the structure and relations of a Bayesian network based on experimental data are given. The lecture ends with an example application of Bayesian networks to gene expression analysis.

¹Based on scribe by Erez Yaffe and Dan Cohen, June 2004 and on scribe by Tal Peled and David Burstein, December 2003

²This lecture is based on presentations of Nir Friedman

8.1.2 A probability primer

Basic notions and notations

In the following lecture we will deal with discrete probabilities defined over a finite *sample space*. The **sample space** Ω is a group containing all of the possible **world states**. A world state defines everything that can be known about the world and the value of any measurable parameter. Any biological experiment is unaware of the current world state and can only measure different parameters resulting from it. Each world state has a given probability to occur and this probability may differ from one world state to another.

A random variable X is a parameter that has a given value for each world state (i.e. it is a function on the group Ω). $P(X = x)$ is the probability for X to be measured as x in a random world state. Thus $P(X = x)$ is the sum of probabilities of the world states in which X will be measured to be x . Many times, the name of the random variable is omitted so $P(x)$ will be used instead of $P(X = x)$.

Given a group of values A for a random variable X , we will note the chance that one of these values will be measure for X by $P(X = A)$ or $P(A)$. It is obvious that if $A = x_1 \dots x_n$ then $P(A) = \sum_{i=1}^n P(x_i)$. Notice also that if X and Y are completely unrelated, i.e. *independent*, then $P(x, y) = P(x) \cdot P(y)$ and $P(A, B) = P(A) \cdot P(B)$. This independence will be noted $I(X; Y)$ or $I(A; B)$ when A is related to X and B to Y .

Joint and conditional distribution

Given two random variables, X and Y , we would like to know their joint distribution $P(X = x, Y = y)$, i.e. the probability that in the same time X is measured to be x and Y is measured to be y . As before, this could be written as $P(x, y)$ even though x and y are values of (possibly) very different parameters. If A and B are groups of possible values for X and Y , respectively, $P(X = A, Y = B)$ or $P(A, B)$ is the probability for X to have a value from A and Y from B . It is obvious that if $A = \{x_1 \dots x_n\}, B = \{y_1 \dots y_m\}$ then $P(A, B) = \sum_{i=1, j=1}^{i=n, j=m} P(x_i, y_j)$.

We would like now to know the probability of measuring $X = x$ given that we have already measure Y to be y . It is obvious that the probability may differ from $P(X = x)$. Think of the probability of someone speaking Hebrew if we already knew that he is from Mars. It should be much smaller than the overall probability for him to speak Hebrew (unless there is something we don't know about Mars...). The above probability is the conditional probability $P(X = x|Y = y)$ or $P(x|y)$. It is logical that the probability of $X = x$ and $Y = y$ should be equal for the probability for " $Y = y$ " and " $X = x$ given that $Y = y$ ". These are

completely unrelated variables (even if X and Y are related) so :

$$P(x, y) = P(y)P(x|y) \quad (8.1)$$

and

$$P(A, B) = P(B)P(A|B) \quad (8.2)$$

We can interchange X and Y and get $P(A, B) = P(A) \cdot P(B|A)$. From these results the important Bayes rule can be inferred :

$$P(A | B) = \frac{P(B | A) \cdot P(A)}{P(B)} \quad (8.3)$$

Also notice that if $I(A;B)$, that is X and Y are independent then

$$P(A | B) = \frac{P(A, B)}{P(B)} = \frac{P(A) \cdot P(B)}{P(B)} = P(A) \quad (8.4)$$

This makes sense, because X doesn't depend on Y so knowing Y doesn't change the probability distribution of X .

Conditional Independence

It is possible that two random variables X and Y are dependent, but given a value of some other random variable Z they become independent. This is the case when X and Y are both dependent on Z but otherwise unrelated to each other. Once we fix the value of Z , there is no remaining information to create dependence between the distributions of X and Y . But without fixing Z we could infer from the value of X the value of Z and using the value of Z to learn about the possible values of Y so they are dependent. This case is called Conditional Independence and will be marked $I(X; Y | Z)$ or $I(A; B | C)$. In this case it is obvious that :

$$P(A, B | C) = P(A | C) \cdot P(B | C) \quad (8.5)$$

Or alternatively:

$$P(A | B, C) = P(A | C), P(B | A, C) = P(B | C) \quad (8.6)$$

Note that $I(A; B) \Rightarrow I(A; B | C)$: if A and B are independent, then they will be independent given any known C . Nevertheless, $I(A; B | C) \nRightarrow I(A; B)$: if A and B are independent given known C , it does *not* mean A and B are independent in any case.

Total probability Theorem

The idea is to break up the sample space according to the value of one random variable Y . Suppose that B is the set composed of the possible values for Y . Let us break B into a number of disjoint sets B_i so that $\bigcup_{i=1}^n B_i = B$. This breaks the sample space accordingly. The total probability theorem states that

$$P(A) = \sum_{i=1}^n P(A, B_i) = \sum_{i=1}^n P(B_i) \cdot P(A | B_i) \quad (8.7)$$

The chain Rule

Using the fact that $P(X_1, \dots, X_n) = P(X_1 | X_2, \dots, X_n) \cdot P(X_2, \dots, X_n)$ we could write $P(X_1, \dots, X_n)$ in the following way :

$$P(X_1, \dots, X_n) = P(X_1 | X_2, \dots, X_n) \cdot P(X_2 | X_3, \dots, X_n) \cdot \dots \cdot P(X_{n-1} | X_n) \cdot P(X_n) \quad (8.8)$$

This equation is called the *chain rule*.

8.1.3 Exploiting Independence Property

Independence of variables can be used to simplify significantly the representation of complex joint distributions. Consider the representation of the joint distribution of n variables, for example n tosses of a coin. Even in this simple binary example, explicit representation of the joint distribution will require 2^n entries, for all possible assignment of heads / tails to the tosses.

Example 8.1: Presenting Distribution of four Fair Tosses

Let $X_i = 1$ denote heads in toss number i . The joint distribution of four tosses is represented in table 8.1.

This is obviously a wasteful representation. If the independence of each toss is exploited, a much more compact representation can be achieved. Such a representation would include only the probabilities of each toss (table 8.2). In this degenerate case the tosses are of a fair coin, so the probability to get heads in each toss is 0.5.

The joint distribution $P(x_1, x_2, x_3, x_4)$ would be computed using this table as $P(X_1 = x_1) \cdot P(X_2 = x_2) \cdot P(X_3 = x_3) \cdot P(X_4 = x_4)$. In most biological applications, the majority of the variables are not completely independent, but exploiting the existing independencies can be used to achieve a rather simple representation of complex dependencies. First a simple example will be presented to explain the principles of this method.

X_1	X_2	X_3	X_4	$P(X_1, X_2, X_3, X_4)$
0	0	0	0	0.0625
0	0	0	1	0.0625
0	0	1	0	0.0625
0	0	1	1	0.0625
0	1	0	0	0.0625
0	1	0	1	0.0625
0	1	1	1	0.0625
1	0	0	0	0.0625
1	0	0	1	0.0625
1	0	1	0	0.0625
1	0	1	1	0.0625
1	1	0	0	0.0625
1	1	0	1	0.0625
1	1	1	1	0.0625

Table 8.1: Joint distribution representation of four fair tosses

	$X_i = 0$	$X_i = 1$
X_1	0.5	0.5
X_2	0.5	0.5
X_3	0.5	0.5
X_4	0.5	0.5

Table 8.2: Independent distribution representation of four fair tosses

Example 8.2: Two Variables Conditional Probability Distribution (CPD)

Consider the problem of determining whether a suspect is guilty. The suspect is standing trial, and the judge finds him either guilty or innocent. The probability space is the joint distribution over these two variables: whether the suspect is **g**uilty, G , and whether the **j**udge will find him guilty, J . Assuming these two variables are binary (0 - not guilty; 1 - guilty), the joint distribution has four entries. Denote the case that $G = 0$ by g^0 , the probability the suspect is not guilty, and in the same manner we use g^1 , j^0 , j^1 . A possible joint distribution of G and J is described in table 8.3.

This is the *joint distribution representation* which presents all the different assignments of valid values to the variables: The probability that the suspect is innocent and the judge finds him innocent is 0.38, the probability that he is innocent, and yet found guilty by the

G	J	$P(G, J)$
g^0	j^0	0.38
g^0	j^1	0.02
g^1	j^0	0.06
g^1	j^1	0.54

Table 8.3: Joint distribution of two variables (the suspect-judge problem)

judge is 0.02, the probability that the suspect is guilty but found innocent is 0.06, etc. . .

The same joint distribution can be alternatively displayed in a somewhat more natural manner. Instead of specifying the various joint entries, the product rule (equation 8.2) is used to get $P(J, G) = P(J | G) \cdot P(G)$, so $P(G)$ and $P(J | G)$ are to be represented. This will require the use of two tables, one representing the distribution of G , and the other representing the *conditional probability distribution (CPD)* of J given G . This representation is denoted as a *factorial representation* and is presented in table 8.4.

G	$P(G)$
g^0	0.4
g^1	0.6

$P(J G)$	j^0	j^1
g^0	0.95	0.05
g^1	0.1	0.9

Table 8.4: Factorial representation of two variable distribution (suspect-judge problem)

The distribution $P(J | G)$ represents the probability for the judge to find the suspect guilty in the two possible alternatives: In the case of g^0 - suspect not guilty and in case of g^1 - suspect is guilty. It is easy to see in this representation the false negative probability: $P(j^0 | g^1)$ - the judge found the suspect innocent although he was guilty (10% in this example), and the false positive probability: $P(j^0 | g^0)$ - the probability of an innocent suspect to be found guilty (5%).

It's also easy to switch between the tables using the product rule (equation 8.2), for example in the joint distribution representation appears $P(g^0 j^0) = 0.38$ (suspect found innocent rightfully). This can be achieved from the factorial representation using the product rule $P(J, G) = P(J | G) \cdot P(G)$, $P(g^0 j^0) = P(j^0 | g^0) \cdot P(g^0) = 0.95 \cdot 0.4 = 0.38$. Switching from joint distribution representation to factorial is done simply by summing the relevant probabilities: $P(g^0) = P(g^0 j^0) + P(g^0 j^1) = 0.38 + 0.02 = 0.4$.

It is instructive to consider the number of independent parameters needed in each of the

representations. The joint distribution representations requires three independent parameters, as there are four entries, and their sum is known to be 1, so the fourth completes to 1. The factorial representation consists of three probability distribution:

1. $P(G)$, so $P(g^0) + P(g^1) = 1$.
2. $P(J | g^0)$, so $P(j^0 | g^0) + P(j^1 | g^0) = 1$.
3. $P(J | g^1)$, so $P(j^0 | g^1) + P(j^1 | g^1) = 1$.

Therefore, only three parameters are required. Note that in this example the two representations use the same number of parameters, but as will be seen further, the factorial method enables a more compact representation.

Bayesian networks will be formally defined further on, in Section 8.2, yet it will be useful to consider how a Bayesian network of this example will be represented. A Bayesian network in this case will have two nodes, representing the two variables G and J , and a single edge from G to J , representing J depends on G (the *direction of dependence* in this model is from J to G). Figure 8.1 shows the Bayesian network in this simple case:

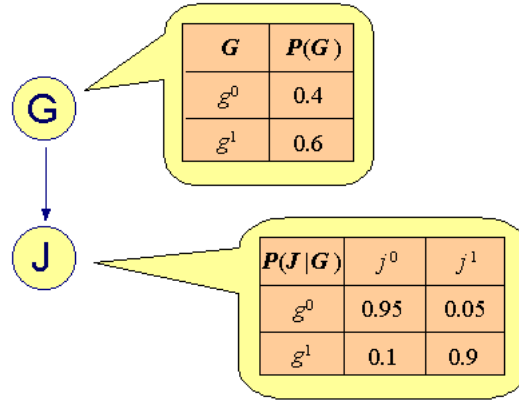


Figure 8.1: Bayesian network for the suspect-judge example. Each variable in the network is associated with a local probability distribution

Example 8.3: Three Variables Conditional Probability Distribution (CPD)

Now a more complicated setting will be considered, where the suspect is being tested on a polygraph during the investigation. As polygraph results are not adequately reliable, they are not admissible in court, and therefore the judge does *not* see the results. This means his verdict is independent of the polygraph's results. Now there are three variable to consider: G , J and the outcome of the polygraph, L . Assuming the polygraph results are also binary,

the joint distribution has eight entries. Table 8.5 presents a possible joint distribution of G , J and L .

G	J	L	$P(G, J)$
g^0	j^0	l^0	0.304
g^0	j^0	l^1	0.076
g^0	j^1	l^0	0.016
g^0	j^1	l^1	0.004
g^1	j^0	l^0	0.006
g^1	j^0	l^1	0.054
g^1	j^1	l^0	0.054
g^1	j^1	l^1	0.486

Table 8.5: Joint distribution of three variables (the suspect-judge-polygraph problem)

The result of the polygraph and the judge's verdict clearly depends on whether the suspect is guilty or not. The polygraph outcome and the verdict are also dependent: if it is given that one of them found the suspect guilty, it increases the probability that the other will also find him guilty.

However, there is a *conditional* independence: If it is known whether the suspect is guilty, the outcome of the polygraph does no longer give information about the verdict, and vice versa, assume it is known the suspect is not guilty, the fact that the judge found him guilty, for example, does not increase the probability the polygraph will show the same. Denoted by l^1 the polygraph outcome which states that the suspect is guilty and l^0 otherwise. Then, $P(l^1 | g^1, h^1) = P(l^1 | g^1)$. More generally it can be assumed that $I(J; L | G)$, this means if G is given, L and J are independent. Note that this assumption is based on the fact that the judge and the polygraph take into considerations completely different aspects of the case.

The conditional independence allows to provide a compact representation of the joint distribution, as seen before. Based on the product rule (8.2), the joint distribution is:

$$P(G, J, L) = P(J, L | G) \cdot P(G)$$

The conditional probability assumed implies that:

$$P(J, L | G) = P(J | G) \cdot P(L | G)$$

Hence,

$$P(G, J, L) = P(J | G) \cdot P(L | G) \cdot P(G) \tag{8.9}$$

Equation 8.9 leads to the desired factorial representation. Therefore, in order to specify fully the joint distribution of this case, we have to know $P(G)$, $P(J | G)$ and $P(L | G)$. Note

that $P(G)$ and $P(G \mid J)$ distributions stays the same as in example 8.2 (table 8.4). This means that expanding the previous example to include also L required only the distribution of $P(L \mid G)$ to be added, this distribution is displayed in table 8.6. That is in contrast to the eight entries joint distribution representation, in which the whole distribution has to be recalculated.

$P(L \mid G)$	l^0	l^1
g^0	0.8	0.2
g^1	0.1	0.9

Table 8.6: The conditional $P(L \mid G)$ distribution which together with table 8.4 ,matches the joint distribution described in table 8.5.

Together these conditional probability distributions fully specify the joint distribution. For example: $P(g^0, j^0, l^1) = P(g^0) \cdot P(j^0 \mid g^0) \cdot P(l^1 \mid g^0) = 0.4 \cdot 0.95 \cdot 0.2 = 0.076$. Unlike the previous example, this one demonstrates a situation in which the factorial representation is more compact: only five independent parameters suffice to fully specify the whole distribution with conditional dependencies (for example $P(g^0)$, $P(j^0 \mid g^0)$, $P(j^0 \mid g^1)$, $P(l^0 \mid g^0)$, $P(l^0 \mid g^1)$, since the other only completes to 1), whilst seven parameters are required using the joint representation (the eighth parameter completes to 1). Note that we succeed to reduce the number of parameters due to conditional independence between the variables. Moreover, the factorial representation has another important property, which was mentioned before, *modularity*. When adding a new variable to a model, L in this example, only the local probability model for L needs to be added, while the local probabilities of G and J are reused. In joint distribution, this is not the case, and all the probabilities have to be recalculated, as the distribution changes entirely.

The Bayesian network corresponding to this scenario (figure 8.2) contains an additional node in its network, representing L , and an edge from G to L , denoting the direction of dependence: L depends on G . Note that given the value of the parent G , both L and J are independent.

8.2 Bayesian Networks

8.2.1 Representing Distributions with Bayesian Networks

Suppose we are given a set of assertions and a variety of ways in which they support each other. Each assertion establishes a value for an attribute and is of the form $(X_i = x_i)$, that is, "Variable X_i has value x_i ". The variables are X_1, \dots, X_n . We would know everything we need to know about the world described by these assertions if we had the joint probability

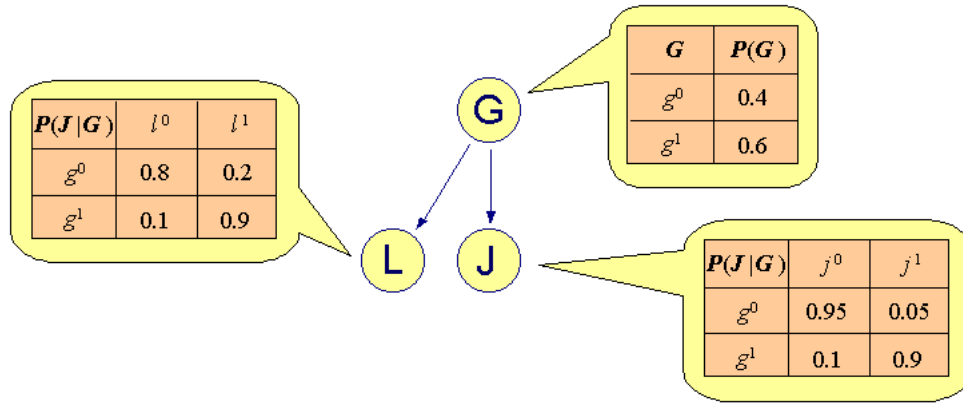


Figure 8.2: Bayesian network for the suspect-judge-polygraph example

$P(X_1, \dots, X_n)$. From this probability function we could compute any other probability such as $P(X_2)$ or $P(X_2 \mid X_3, X_5)$. Unfortunately, even when assuming for simplicity that the variables are binomial, the representation complexity of $P(X_1, \dots, X_n)$ is, as we've seen, 2^n , which is impractical even for small value of n .

Bayesian Networks simplify this problem by taking advantage of existing causal connections between assertions, and of assumptions about conditional independence. A *Bayesian network* is a representation of a joint probability distribution. This representation consists of two components:

- The first component, G , is a *directed acyclic graph* (DAG) whose nodes correspond to the random variables X_1, \dots, X_n , and its edges correspond to dependencies and their directions. This component is known as the *Qualitative part*.
- The second component describes the *local probability model, the conditional probability distribution (CPD)* for each variable, given its parents in G . Let X_i be a variable and $\mathbf{Pa}(X_i)$ its parents in G , the CPD of X_i is the distribution of $P(X_i \mid \mathbf{Pa}(X_i))$. This part is the *Quantitative part*.

Together, these two components specify a unique distribution over X_1, \dots, X_n . The graph G represents conditional independence assumptions that allow the joint distribution to be decomposed, economizing on the number of parameters. Note that the graph G encodes the *Markov Assumption*: *Each variable X_i is independent of its non-descendants, given its parents in G* (it is still dependent on its descendants). It is a natural assumption for many causal processes. If the parents of an event, meaning the events which directly caused it, are known, it is independent of the events which effect its parents, or of any other event, except the events which it is the cause for, or the cause for one of their ancestors.

By applying the chain rule of probabilities and properties of conditional independencies, any joint distribution can be decomposed to the *product form* according to the Markov assumption on it's Bayesian network. Suppose w.l.o.g (without limitation of globality) X_1, \dots, X_n are arranged in reversed topological order, i.e. if X_j is a descendant of X_i in the network then $j < i$. According to the chain rule (equation 8.8):

$$P(X_1, \dots, X_n) = P(X_1 \mid X_2, X_3, \dots, X_n) \cdot P(X_2 \mid X_3, X_4, \dots, X_n) \cdot \dots \cdot P(X_{n-1} \mid X_n) \cdot P(X_n)$$

Since each X_i does not depend on any of its non-descendants given its parents, and any variable X_j that is a descendant of X_i has a lower index ($j < i$) due to the reversed topological order, then $P(X_i \mid X_{i+1}, \dots, X_n) = P(X_i \mid \mathbf{Pa}(X_i))$, where $\mathbf{Pa}(X_i)$ is the set of parents of X_i in graph G . Thus, we get the chain rule for Bayesian networks:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid \mathbf{Pa}(X_i)), \quad (8.10)$$

Example 8.4: A Simple Bayesian Network with Five Variables

In Figure 8.3 we can see a simple example of a Bayesian network structure. This network describes the connections between the following events:

- B - There is a **B**urglary.
- A - The **A**larm goes off.
- E - There is an **E**arthquake.
- R - There is a **R**adio report of an earthquake.
- C - Mr. Watson, the neighbor, **C**alls to inform us he heard the alarm.

The alarm is set to detect earthquakes and attempts of burglary and. Naturally, there is a probability of a false alarm to occur, or for some malfunction to cause the alarm not to operate when it should. Mr. Watson is the neighbor, and in case he hears the alarm he should call us to inform the alarm had gone off. He might, of course, not hear it or mistakenly think he had heard it (though it hadn't gone off). In addition the local radio station usually reports any case of an earthquake.

Specifying the joint distribution of these events, requires $2^5 - 1 = 31$ parameters. This is quite a high number for such a simple system. Consider an expert system for monitoring intensive care patients, which measures over 37 different properties. Joint distribution representation of such a system will require at least 2^{37} parameters. This is not feasible. Instead, we will use a Bayesian network to represent such systems (figure 8.3).

Following is the list of independencies of this network:

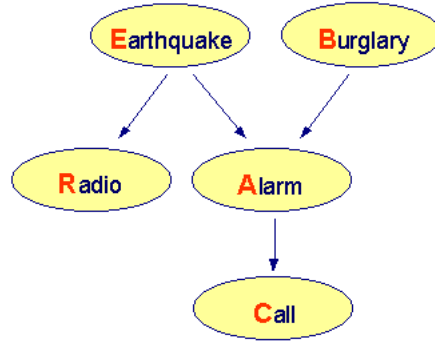


Figure 8.3: Simple Bayesian Network with Five Variables

- $I(E, B)$ - Given the parents of E (none), E is conditionally independent of its non-descendant: B .
- $I(B, (E, R))$ - Given the parents of B (none), B is conditionally independent of its non-descendants: E and R .
- $I(R, (B, A, C) \mid E)$ - Given the parents of R (E), R is conditionally independent of its non-descendants: B , A and C .
- $I(A, R \mid (B, E))$ - Given the parents of A (B, E), A is conditionally independent of its non-descendant: R .
- $I(C, (R, B, E) \mid A)$ - Given the parents of C (A), C is conditionally independent of its non-descendants: R , B and E .

Consider the first independence, $I(E, B)$: E has no parents, so in any case it is independent of B . Indeed there is no dependence between the events of a burglary and an earthquake. The last independence, $I(C, (R, B, E) \mid A)$ means that if the parents of C , which is A , is known then C is independent of R , E and B . This correlates with the rational of the events: if it is known that the alarm broke off (A), then whether the neighbor will or will not call (C) because of it, does not depend on the radio report (R), the earthquake (E) or the burglary (B). This is true despite the fact that the earthquake and the burglary might be the reason for the alarm.

According to the chain rule (equation 8.8), the joint distribution is:

$$P(C, A, R, E, B) = P(C \mid A, R, E, B) \cdot P(A \mid R, E, B) \cdot P(R \mid E, B) \cdot P(E \mid B) \cdot P(B)$$

Which requires 31 parameters. Alternatively, using independencies in the Bayesian network (equation 8.10), the joint distribution of the five events is

$$P(C, A, R, E, B) = P(C \mid A) \cdot P(A \mid B, E) \cdot P(R \mid E) \cdot P(E) \cdot P(B)$$

This distribution requires only 10 independent parameters.

Complexity Analysis

To fully specify a joint distribution, we need to specify the conditional probabilities in the product form. The quantitative part of the Bayesian network describes these conditional distributions, $P(X_i \mid \mathbf{Pa}(X_i))$ for each variable X_i . Suppose the variables have k possible values (If the variables were binomial then $k = 2$). Each one of these conditional distributions contains $k^{|\mathbf{Pa}(X_i)|}$ independent parameters. Let l be the maximum number of possible parents, then for each variable X_i , there are $\leq k^l$ parameters to be stored. Hence for n variables, the representation complexity is $O(n \cdot k^l)$. Note that for a joint distributions with many variables, each of which has only few dependencies, then $l \ll n$ and the representation complexity is much better than the $O(k^n)$ representation complexity in the joint distributions representation.

8.2.2 Representing Markov chains and HMM as Bayesian networks

A Markov chain can be treated as a simple case of a Bayesian network. A Markov chain



Figure 8.4: A Markov chain represented as a Bayesian network

is a triplet Q, p, A , where Q is a set of states, p is the probabilities of the initial state, and A is the state transition probabilities, which contains for each two states $s, t \in Q$ the probability $a_{st} \equiv P(x_i = t \mid x_{i-1} = s)$. $X = (x_1, \dots, x_n)$ is a random process. An equivalent Bayesian network will be a chain of variables X_1, \dots, X_n (figure 8.4), thus X_{i-1} is the parent of X_i . Each variable may attain any value $s \in Q$. The CPD of X_1 will be p . The CPD of variable X_i is $P(X_i \mid X_{i-1}) = a_{st}$, i.e. the probability of observing value s in variable X_i is dependent on the value in its parent X_{i-1} in the Bayesian network. Therefore computing $P(X) = P(x_1) \cdot \prod_{i=2}^L a_{x_{i-1}x_i}$ using Markov chain, is equivalent to the factorial representation of the joint distribution in the Bayesian network: $P(X) = P(X_1 = x_1) \cdot \prod_{i=2}^L (X_i = x_i \mid X_{i-1} = x_{i-1})$.

Hidden Markov model (HMM), can also be represented using Bayesian networks. In addition to the Markov chain, it contains emission probabilities, which will be represented as additional nodes (Y_i) in the Bayesian network (figure 8.5). The emission node will have no descendants, as nothing depends on them. Each variable Y_i may attain any value $y_i \in \Sigma$. The CPD of Y_i will be represented as the emission probability, i.e. $P(Y_i = y_i \mid X_i = x_i) = e_{x_i}(y_i)$. The CPD of each variable X_i remains the same as in the Markov chain model.

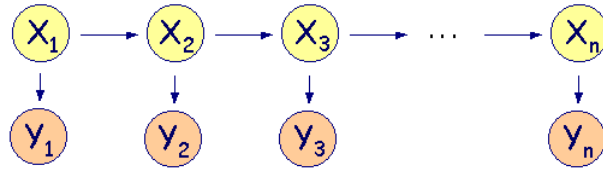


Figure 8.5: A HMM as a Bayesian network

As demonstrated, it's possible to represent HMM using Bayesian network, but this representation is wasteful, since a_{kl} and $e_k(b)$ appear identically n times, once in each node. In the CpG problem, for example, locating CpG islands in a sequence of 2,000 nucleotides, will require an HMM with two states, but the Bayesian network of the same problem will contain 2,000 nodes representing the probabilities of the states in each one of the nucleotides, and another 2,000 for the emissions probabilities. Therefore for random processes with memory of length 1, it will be easier to use HMM rather than Bayesian networks. However, in HMM each state depends only on the state before it, and more complicated dependencies cannot be expressed. As the Bayesian network model generalizes, and represents high-order dependencies, it should be used in cases of complicated dependencies.

8.3 Inference in Bayesian networks

8.3.1 Introduction to Inference

So far it has been shown that Bayesian networks can be used to represent probability distributions in a compact and intuitive manner. As it is so, Bayesian networks should contain information that would answer any query about the distributions represented by them. This section is devoted to the way such queries can be answered, using the data available by the network. Given a Bayesian network for the variables X_1, \dots, X_n , frequent types of queries we are interested in are $P(x_i)$, $P(x_i | X_j = x_j)$, $P(x_1, \dots, x_n | X_j = x_j)$ etc.

Bayesian networks can be used to answer these queries, since the joint distributions can be generated using the chain rule for Bayesian networks (equation 8.10). The joint distribution can then be used to make the necessary calculations for the query. A naïve solution to compute $P(X_i)$ based on the total probability theorem (equation 8.7), is $P(X_i) = \sum_{x_1} \cdots \sum_{x_{i-1}} \cdot \sum_{x_{i+1}} \cdots \sum_{x_n} P(X_1, \dots, X_n)$. Using the joint distribution representation, for binomial variables, the complexity of solving this is exponential ($O(2^n)$). For variables with up to k possible values, it is $O(k^n)$. Avoiding exponential complexity, was one of the main reasons to use Bayesian networks, and in the following section it will be demonstrated how the use of them might reduce the complexity. The inference problem in Bayesian networks is NP-hard. It means that in the *worst case* the best algorithm available

to solve these problems is exponential (assuming $P \neq NP$). In other words - performing the computations on all the joint distribution entries is the best that can be done in worst case. But it turns out that the worst case comes up only very rarely. So, although the complexity of Bayesian networks algorithms in *worst case* is exponential, in practice, a simple method can be used to make the computations in most practical cases quite fast. This technique is called *Variable elimination* algorithm.

8.3.2 Variable Elimination algorithm

The Basic Idea

The idea of eliminating variables will be demonstrated through a basic inference task on a simple Bayesian network. Consider the “chain” network in figure 8.6. Suppose we want to



Figure 8.6: Simple Bayesian network in form of a nodes chain

calculate $P(A = a, D = d)$, the distribution of the joint probability of A and D : Using the total probability equation (8.7), we get:

$$P(a, d) = \sum_b \sum_c P(a, b, c, d)$$

Note that each addition is used to eliminate one variable: Summing for all the probabilities of c on the joint distribution, $P(a, b, c, d)$, will yield $P(a, b, d)$, thus eliminating C . In the same manner summing up all the probabilities of b on $P(a, b, d)$ will eliminate b , yielding the requested probability $P(a, d)$. Using the factorial representation, instead of the joint one, will give us:

$$P(a, d) = \sum_b \sum_c P(a, b, c, d) = \sum_b \sum_c P(d | c)P(c | b)P(b | a)P(a)$$

This representation embodies the power of the Bayesian network, as it uses the conditional probabilities instead of the joint, and these will be used to simplify the additions and reduce the complexity.

First, c will be eliminated. Note that when summing for all $C = c$ on the conditional probabilities, some terms need not to be summed, as their CPDs do not contain the variable C . That means we can omit these terms from the sum on c , thus reducing the complexity, as less calculation has to be performed to accomplish that sum:

$$P(a, d) = \sum_b \sum_c P(d | c)P(c | b) \overbrace{P(b | a)P(a)}^{\text{No } c \text{ in CPDs}} = \sum_b P(b | a)P(a) \sum_c P(d | c)P(c | b)$$

Now, in order to complete the elimination of c , we will use the chain rule (equation 8.8), applying $\sum_c P(d | c)P(c | b) = P(d | b)$:

$$P(a, d) = \sum_b P(b | a)P(a) \overbrace{\sum_c P(d | c)P(c | b)}^{P(d|b)} = \sum_b P(b | a)P(a)P(d | b)$$

c was eliminated from this equation, and the factor of $P(d | b)$ was calculated. $P(d | b)$ is an *intermediate factor* of the algorithm. The next step will be to eliminate B . As before, this will be done by summing on every $B = b$, after excluding the terms that do not include b (in this case $P(a)$) from the addition.

$$P(a, d) = \sum_b P(b | a) \overbrace{P(a)}^{\leftarrow} P(d | b) == P(a) \overbrace{\sum_b P(d | b)P(b | a)}^{P(d|a)} = P(d | a)P(a)$$

b was eliminated from this equation, and the intermediate factor $P(d | a)$ was calculated. Next, in order to compute $P(a, d)$, we only need to multiply $P(a)$ by the intermediate factor $P(d | a)$, which was already calculated.

Generalization

The technique presented can be generalized to solve any specific distribution of a variable, or joint distribution of a subset of variables: Let $P(x_1)$ be the distribution we're looking for in a network of n variables. First the query will be written in the following form:

$$P(x_1) = \sum_{x_n} \cdots \sum_{x_3} \sum_{x_2} \prod_i P(x_i | pa_i)$$

Where pa_i are the parents of x_i in the Bayesian network. In the last example, $\prod_i P(x_i | pa_i) = P(d | c)P(c | b)P(b | a)P(a)$, After having the query in that form, perform iteratively:

- Move all irrelevant terms outside of innermost sum.
- Perform innermost sum, getting a new intermediate factor.
- Insert the intermediate factor into the product

So far we've seen how to answer a query about the distribution of a variable, or the joint distribution of few variables. Consider a conditional probability query. Let X_1 be the query variable given $X_k = x_k$, so the distribution of $P(X_1 | x_k)$ is requested.

The query will be managed in a similar way as in the previous case: First $P(X_1, x_k)$ and $P(x_k)$ will be calculated, as follows. Note that for $P(X_1, x_k)$ we sum neither on X_1 nor on x_k .

$$P(X_1, x_k) = \sum_{x_n} \cdots \sum_{x_{k+1}} \sum_{x_{k-1}} \cdots \sum_{x_2} \prod_i P(x_i | pa_i)$$

$$P(X_k = x_k) = \sum_{x_n} \cdots \sum_{x_{k+1}} \sum_{x_{k-1}} \cdots \sum_{x_1} \prod_i P(x_i | pa_i)$$

Next $P(X_1|X_k)$ is found, using the product rule (equation 8.2): $P(X_1 | x_k) = \frac{P(X_1, x_k)}{P(X_k = x_k)}$

Complexity Analysis

Let's observe the complexity of solving the queries presented. As mention in section 8.3.1, solving such queries in a Bayesian network is NP-hard. Therefore, in the worst case, time complexity is expected to be exponential. But most practical cases do not yield such networks. In the following calculations the variables will be treated at first as binomial, to simplify the calculations. In order to find the distribution of X_1 , the following expression needs to be calculated.

$$P(x_1) = \sum_{x_n} \cdots \sum_{x_3} \sum_{x_2} P(x_2, x_3, \dots, x_n)$$

In the naïve case, all the additions will be performed on the joint distribution. Since we sum over all the possible combination of x_1, \dots, x_n , this will cost $O(2^n)$. In the multinomial case, where X_i has k possible values, the complexity will be $O(k^n)$.

We've seen that time complexity might be reduced, if the factorial representation is used and each sum is done only on the relevant factors. In this method we start with the following expression:

$$P(x_1) = \sum_{x_n} \cdots \sum_{x_3} \sum_{x_2} \prod_i P(x_i | pa_i)$$

And then perform $O(n)$ additions, each on its relevant factors. Let's observe the complexity of the X_i -th summation: We sum only on the relevant factors of the product, let r_i be the number of variables in the intermediate factor generated in summation of x_i . For example, summing $\sum_c P(A | c) \cdot P(c | b, d)$ will generate the intermediate factor $P(A | b, d)$ with $r = 3$. In the intermediate factor there are $O(2^{r_i})$ possible combinations, and each one is calculated by multiplying $O(r_i)$ factors. Therefore summing X_i will cost $O(r_i \cdot 2^{r_i})$. Let $r = \max_{i=1}^n \{r_i\}$. Summing $O(n)$ variables will cost $O(n \cdot r \cdot 2^r)$. In the multinomial case, the complexity will be $O(n \cdot r \cdot k^r)$. Note that this is significantly smaller than $O(k^n)$ only if each variable has a small number of parents in the network (few parents means few dependencies, which will give factors with few variables, hence a low r).

Finally note that the number of factors, and hence the complexity of the elimination, strongly depend on the order of elimination. Even in the simple elimination example on the

chain network described in figure (8.6), containing only the four variables A, B, C and D this fact is manifested. consider the following order of elimination:

$$\begin{aligned}
 P(d) &= \sum_c \sum_b \sum_a P(a, b, c, d) = \sum_c \sum_b \sum_a P(d | c) P(c | b) P(b | a) P(a) \\
 &= \sum_c \sum_b P(d | c) P(c | b) \overbrace{\sum_a P(b | a) P(a)}^{P(b)} = \sum_c \sum_b P(d | c) P(c | b) P(b) \\
 &= \sum_c P(d | c) \overbrace{\sum_b P(c | b) P(b)}^{P(c)} = \sum_c P(d | c) P(c) = P(d)
 \end{aligned}$$

Each one of the products being summed contains 1 factor at most: Just one variable except the one that we sum upon. Alternatively, consider the following order of elimination:

$$\begin{aligned}
 P(d) &= \sum_a \sum_b \sum_c P(a, b, c, d) = \sum_a \sum_b \sum_c P(d | c) P(c | b) P(b | a) P(a) \\
 &= \sum_a \sum_b \overbrace{\sum_c P(d | c) P(c | b)}^{P(d|b)} P(b | a) P(a) = \sum_a \sum_b P(d | b) P(b | a) P(a) \\
 &= \sum_a \overbrace{\sum_b P(d | b) P(b | a)}^{P(d|a)} P(a) = \sum_a P(d | a) P(a) = P(d)
 \end{aligned}$$

In this elimination all the products except the last contains 2 factors, hence the complexity of this elimination is greater, since the order of elimination wasn't optimal. Choosing an efficient order of elimination can be done using a *clique tree algorithm*, which will not be discussed here.

To summarize, the complexity of inferring distributions from the Bayesian network representation obeys the following rules :

- The naive method of inference is exponential in the number of variables. In general, the inference problem is NP-hard.
- Using variable elimination method is exponential in the size of the largest summation factor and polynomial in the number of variables.
- Complexity of variable elimination depends on order of elimination.

8.4 Learning Bayesian Networks

8.4.1 Introduction [8]

The amount of available information is growing rapidly. However, knowledge acquisition is an expensive process. Learning allows us to construct models from raw data, which can then be used for many tasks. Bayesian networks in particular, which convey conditional independencies, enable us to capture the structure of many real-world distributions.

8.4.2 The Learning Problem

The problem of learning a Bayesian network can be stated as follows:

Let m be the number of samples and n the number of variables. Given a training set $D = (X_1, \dots, X_m)$, where $X_i = (x_{i1}, \dots, x_{in})$, and prior information, find a network that *best matches* D . The problem of learning a Bayesian network can be learning the CPD (the parameters) given a certain structure, as in figure 8.7, or learning both the distributions and the graph's structure (the dependencies), as in figure 8.8. In addition, the learning problem can be divided into two other cases: complete data and incomplete data. In complete data, all parameter values are known, whereas in the case of incomplete data, some of the values in the vectors $X_i = (x_{i1}, \dots, x_{in})$ are missing.

The means of handling these aforementioned cases, are described in table 8.7.

	Known Structure	Unknown Structure
Complete Data- All instances of the variables are known	Statistical parametric estimation (closed-form equations)	Discrete optimization over structures (discrete search)
Incomplete Data- Not all instances of the variables are known	Parametric optimization (EM, gradient descent...)	Combined (structural EM, mixture models...)

Table 8.7: Means of learning the network in different cases

We will focus on complete data for the rest of the talk. We start with the case of a known structure, whose parameters we wish to learn and estimate.

8.4.3 The Likelihood Method

One way of finding an estimator for a parameter is the likelihood method. According to this method, given a sequence of samples, we look for the parameter's value that seems most

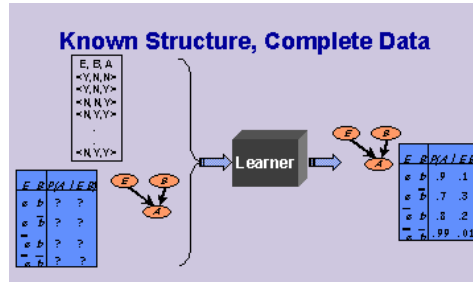


Figure 8.7: Learning from a known structure and complete data

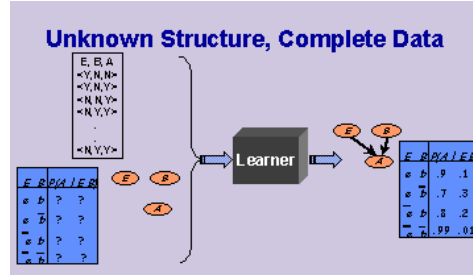


Figure 8.8: Learning from an unknown structure and complete data

”likely” to have caused this certain sequence of samples (i.e. gives our sample the highest probability compared to all other possible values of the parameter). In other words, we try to deduce the value of the parameter **after** a certain sample has been received, namely the value of the parameter that best explains the given sample results.

Likelihood Definitions

- **The Likelihood Function** [17, 13]

Let X_1, X_2, \dots, X_m be samples from a population with a certain distribution with an unknown parameter θ and let $X_i = x[i]$ be the training data- the sample results obtained, each result independent of the others. The likelihood function is defined as:

$$L(\theta : D) = P(D | \theta) = P(X_1 = x[1], X_2 = x[2], \dots, X_m = x[m] | \theta) = \prod_{i=1}^m P(x[i] | \theta)$$

where θ gets all the possible values of the parameter.

- **Maximum Likelihood Estimator**

Given a sample result, $\hat{\theta}$ is the maximum likelihood estimator (MLE) for the parameter θ if it holds that $\forall \theta. L(\hat{\theta}) \geq L(\theta)$. In other words, $\hat{\theta}$ is the value of the parameter θ which maximizes the likelihood function.

8.4.4 Example: Binomial Experiment

A coin can land on one of two positions: H or T, each with an unknown probability. We denote by θ the unknown probability $P(H)$.

Estimation Task: given a sequence of samples $x[1], x[2], \dots, x[m]$ we want to estimate the probability $P(H) = \theta$ and $P(T) = 1 - \theta$.

Suppose we performed six independent Bernolli experiments and received the following sample results: H, H, H, H, T, T. We denote by X_i the result of the i -th coin toss. The likelihood function for the obtained results is:

$$\begin{aligned} L(\theta : D) &= P(D \mid \theta) = P(X_1 = H, X_2 = H, X_3 = H, X_4 = H, X_5 = T, X_6 = T \mid \theta) \stackrel{(*)}{=} \\ &= P_\theta(X_1 = H)P_\theta(X_2 = H)P_\theta(X_3 = H)P_\theta(X_4 = H)P_\theta(X_5 = T)P_\theta(X_6 = T) = \\ &= \theta^4(1 - \theta)^2 \end{aligned}$$

where $(*)$ is because of independent coin tosses.

Note the following:

- For a known θ , we would have an exact numerical value denoting the probability to get the training data results.
- The coefficient $\binom{n}{k}$ in the binomial distribution is omitted here, since we have the exact sample sequence and for each result in the sequence we are interested only in the probability of obtaining that certain result.
In addition, even if we write the coefficient, it will eventually "disappear" from the equation, as explained in the steps to come.

We can formulate the following general case:

Let N_H be the number of occurrences of H and N_T be the number of occurrences of T in the training set, so that $N_H + N_T = m$.

$$L(\theta : D) = P(D \mid \theta) = \theta_H^{N_H} \cdot \theta_T^{N_T}$$

where θ_H is the probability of the outcome H and θ_T is the probability of the outcome T.

As described earlier, the MLE principle is to choose the value of the parameter that maximizes the likelihood function. Applying the MLE principle means we have to find the maximum point of the likelihood function. In order to do this, we will first switch to the log (or \ln) of the likelihood function. This is done in order to make the derivative calculations easier (instead of a complicated derivative of a long product term the log turns the product into summation, thus making the calculation of the derivative simpler). The logarithm is a monotonically increasing function and therefore the maximum point of the function $L(\theta)$ is the same as the maximum point of the function $\log L(\theta)$.

Returning to finding the maximum point of our likelihood function, we get:

$$\ln L(\theta) = 4 \ln \theta + 2 \ln(1 - \theta)$$

$$\frac{d \ln L(\theta)}{d\theta} = \frac{4}{\theta} - \frac{2}{1 - \theta}$$

$$\frac{4}{\theta} - \frac{2}{1 - \theta} = 0$$

$$\hat{\theta} = 2/3$$

or in the general case:

$$\hat{\theta} = \frac{N_H}{N_H + N_T}$$

We can now see that had we written the coefficient, it would have been omitted anyway, after applying the log and finding the derivative.

Using the above general equation, we can easily find the MLE estimation for other cases. For example, given $(N_H, N_T) = (3, 2)$, the MLE estimation is $\frac{3}{5} = 0.6$.

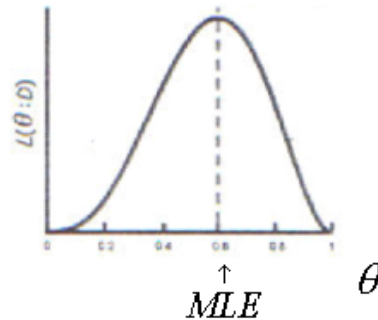


Figure 8.9: The maximum likelihood estimation

8.4.5 Learning Parameters for Bayesian Networks

When learning parameters for Bayesian networks, the training data has a form as in figure 8.10.

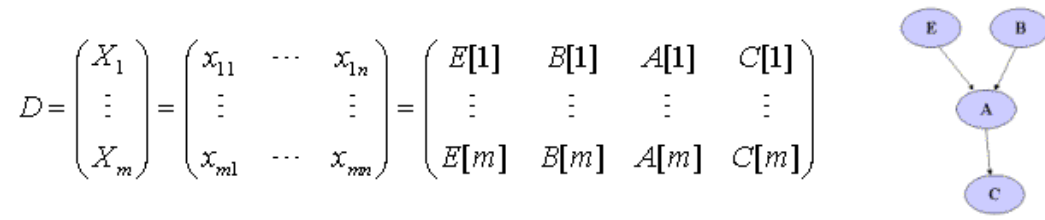


Figure 8.10: The training data, including the network structure

Since we assume i.i.d (independent identically distributed) samples, the likelihood function is:

$$L(\theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \theta) \stackrel{(*)}{=} \underbrace{\prod_m \begin{pmatrix} P(E[m] : \theta) \\ P(B[m] : \theta) \\ P(A[m] | B[m], E[m] : \theta) \\ P(C[m] | A[m] : \theta) \end{pmatrix}}_{\text{the factorial representation}}$$

where (*) is due to switching to the factorial representation, by definition of the network in figure 8.10. Rewriting the terms (looking at the columns of the training data instead of the rows) we get:

$$L(\theta : D) = \prod_m P(E[m] : \theta) \prod_m P(B[m] : \theta) \prod_m P(A[m] | B[m], E[m] : \theta) \prod_m P(C[m] | A[m] : \theta)$$

As can be seen from the calculation above, we have gone from the product of the probabilities of all nodes for each m , to the product of the likelihood of each node (given its parents, if it has any).

Generalizing for any Bayesian network we get:

$$\begin{aligned} L(\theta : D) &= \prod_m P(x_1[m], \dots, x_n[m] : \theta) \underbrace{=}_{\text{network factorization}} \\ &= \prod_i \prod_m P(x_i[m] | Pa_i[m] : \theta_i) \\ &= \prod_i L_i(\theta_i : D) \end{aligned}$$

where Pa_i is the set of node X_i 's parents.

We can see here the **decomposition principle**: the likelihood decomposes according to the structure of the network meaning that we can calculate each node's log-likelihood separately

and sum them all to get the log-likelihood of the entire network. The decomposition enables us to have independent estimation problems. In other words, since the parameters for each variable are not dependent, they can be estimated independently of each other. As a result, we can use the MLE formula we already described and get $\hat{\theta}_{x_i|pa_i} = \frac{N_{x_i|pa_i}}{N_{x_i=0|pa_i} + N_{x_i=1|pa_i}}$ where pa_i is a certain combination of the parents (for example, see figure 8.11).

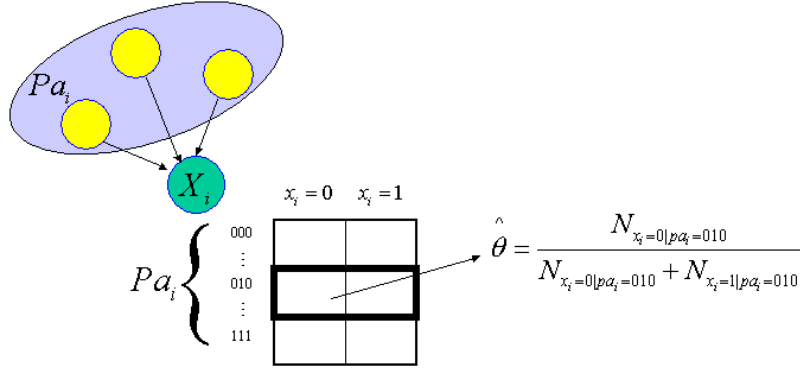


Figure 8.11: Demonstrating the decomposability principle- the variable X_i has three parents, which are the members of the set Pa_i . Each row of the variable's CPD is a different combination of the parents' values and the columns are the different values X_i can receive. Each row is X_i 's distribution given the parents' specific combination (the probabilities in each row sum up to 1) and the parameter is estimated according to the formula for $\hat{\theta}$. In this figure, we wish to estimate the value of $\theta = P(x_i = 0)$ for a given parental combination 010.

8.4.6 Likelihood Function: Multinomials

Up to now we dealt with binomial variables that had either the value 0 or the value 1 ("failure" or "success"). We now make the transition to the multinomial case.

Suppose our variable X can have the values $1, 2, \dots, k$. We would like to learn the parameters $\theta_1, \dots, \theta_k$ (for example each θ_i corresponds to the probability of getting a specific number on a dice). We denote by N_1, \dots, N_k the number of times each outcome is observed. We get the following likelihood function:

$$L(\theta : D) = \prod_{j=1}^k \theta_j^{N_j}$$

where θ_j is the probability of the outcome j and N_j is the number of times the outcome j is observed in D .

As in the binomial case, we have once again omitted the coefficient (this time $\frac{N!}{N_1!N_2!\dots,N_k!}$, $N = N_1 + \dots + N_k$) for it "falls off" after the log likelihood function is derived.

In the same way as before, we get that the MLE is

$$\hat{\theta}_j = \frac{N_j}{\sum_{l=1}^k N_l}$$

In Bayesian networks, when we assume that $P(x_i \mid pa_i)$ is multinomial, we get further decomposition. From equation 8.11 we get:

$$L(\theta_i : D) = \prod_m P(x_i[m] \mid Pa_i[m] : \theta_i)$$

We now write the above product expression in a slightly different way, by grouping products according to a certain combination of the parents' values (for example, in the network depicted in figure 8.10, a combination of A 's parents E and B may be $E = 1$ and $B = 2$). We once again denote by pa_i each combination of the parents and go over all such combinations. Now, for each pa_i we multiply only observations where $Pa_i[m] = pa_i$. Thus,

$$= \prod_{pa_i} \prod_{m: Pa_i[m]=pa_i} P(x_i[m] \mid pa_i : \theta_i)$$

Looking at the second product expression, we now "group" it even further- for each combination of the parents, we take a certain value of x_i and raise its probability to the power of that certain value's number of appearances:

$$= \prod_{pa_i} \prod_{x_i} P(x_i \mid pa_i : \theta_i)^{N(x_i, pa_i)} = \prod_{pa_i} \prod_{x_i} \theta_{x_i|pa_i}^{N(x_i, pa_i)}$$

For each combination pa_i of the parents of x_i we get an independent likelihood function $\prod_{x_i} \theta_{x_i|pa_i}^{N(x_i, pa_i)}$. After performing log on the last expression and finding its derivative we get that the MLE is:

$$\hat{\theta}_{x_i|pa_i} = \frac{N(x_i, pa_i)}{N(pa_i)}$$

This is essentially as in figure 8.11, but for the multinomial case.

8.4.7 Bayesian Learning

Unlike the MLE approach, the Bayesian approach doesn't estimate the most likely θ , but rather takes into account all possible θ s with their probabilities. We can represent our uncertainty about the sampling process using a Bayesian network, as seen in figure 8.12.

All the data are dependent on the certain θ that "created" them. As seen in the Markov

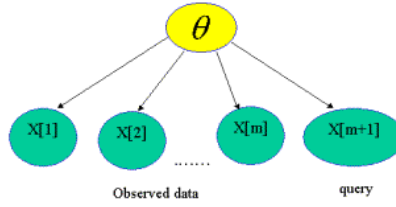


Figure 8.12: Representing Bayesian prediction as inference in Bayesian network

assumption, given θ (the parent) the values of the different $x[i]$ are independent. The Bayesian prediction is inference in this network:

$$P(x[m+1] | x[1], \dots, x[m]) \stackrel{(*)}{=} \int_0^1 P(x[m+1] | \theta, x[1], \dots, x[m]) P(\theta | x[1], \dots, x[m]) d\theta$$

where (*) is because of the conditional probability and total probability theorem for a continuous variable θ ($0 \leq \theta \leq 1$). Since the first probability term in the integral is for a given θ , $x[m+1]$ is independent of its non-descendants and we can therefore write:

$$= \int_0^1 P(x[m+1] | \theta) P(\theta | x[1], \dots, x[m]) d\theta \quad (8.11)$$

Now let us look at the second term of the integral.

$$\underbrace{P(\theta | x[1], \dots, x[m])}_{\text{posterior}} \stackrel{\text{Bayes rule}}{=} \frac{\overbrace{P(x[1], \dots, x[m] | \theta)}^{\text{likelihood}} \overbrace{P(\theta)}^{\text{prior}}}{\underbrace{P(x[1], \dots, x[m])}_{\text{probability of data}}} \quad (8.12)$$

The **prior** is some earlier information we have regarding the random variable θ . Using the total probability theorem, we may rewrite the denominator and get:

$$P(x[1], \dots, x[m]) = \int_0^1 P(x[1], \dots, x[m] | \theta) P(\theta) d\theta$$

and therefore the expression on the right hand side of equation 8.12 can be written as

$$\frac{P(x[1], \dots, x[m] | \theta) P(\theta)}{\int_0^1 P(x[1], \dots, x[m] | \theta) P(\theta) d\theta}$$

Note that $P(x[1], \dots, x[m])$ does not depend on θ (it goes over all values of θ in its equivalent integral expression). As a result, it behaves as a constant in formula 8.11 and can be taken out of the integral there (see the example in the following section, 8.4.8).

8.4.8 Binomial Example: MLE vs. Bayesian Learning

We return to our coin toss example. Recall that $P(H) = \theta$. This time, however, we also have prior information that the distribution of θ is uniform in $[0, 1]$. Therefore, $P(\theta) = f(\theta) = 1$ (no substantial information about θ). In addition, the data is $(N_H, N_T) = (4, 1)$ (there are four heads and one tail).

According to the MLE general case found earlier, the MLE for $P(X = H)$ is $\frac{N_H}{N_H + N_T} = \frac{4}{5} = 0.8$.

The Bayesian prediction is:

$$\begin{aligned}
 P(x[m+1] = H \mid D) &= \int_0^1 P(x[m+1] = H \mid \theta) P(\theta \mid x[1], \dots, x[m]) d\theta = \\
 &= \int_0^1 P(x[m+1] = H \mid \theta) \frac{P(\theta) P(x[1], \dots, x[m] \mid \theta)}{P(x[1], \dots, x[m])} d\theta = \\
 &= \int_0^1 \frac{P(x[m+1] = H \mid \theta) P(\theta) P(x[1], \dots, x[m] \mid \theta) d\theta}{\int_0^1 P(\theta) P(x[1], \dots, x[m] \mid \theta) d\theta} = \\
 &= \frac{\int_0^1 \theta \cdot 1 \cdot \theta^{N_H} (1 - \theta)^{N_T} d\theta}{\int_0^1 1 \cdot \theta^{N_H} (1 - \theta)^{N_T} d\theta} = \frac{\int_0^1 \theta \cdot \theta^4 \cdot (1 - \theta) d\theta}{\int_0^1 \theta^4 \cdot (1 - \theta) d\theta} = \\
 &= \frac{\int_0^1 (\theta^5 - \theta^6) d\theta}{\int_0^1 (\theta^4 - \theta^5) d\theta} = \frac{\frac{1}{6} - \frac{1}{7}}{\frac{1}{5} - \frac{1}{6}} = \frac{5}{7} = 0.7142
 \end{aligned}$$

Since our earlier knowledge was θ 's uniform distribution, we were not "tempted" to believe, after only 5 tosses, that $P(H)$ is as extreme as 0.8, but rather slightly lower, toward the uniform 0.5. We can see that the prior has given us a more "balanced" result. It seems as though we have added two more tosses, one of which was H, namely the prior has given us a more established result, based on our earlier knowledge.

8.4.9 Summary: MLE vs. Bayesian Learning

The differences between the MLE approach and the Bayesian approach are summarized in table 8.8.

8.4.10 Dirichlet Priors

In example 8.4.8, we saw an example of a uniform prior in $[0, 1]$, and thus $P(\theta) = f(\theta) = 1$. Let us now look at the following density function (recall that θ is a random variable), called

MLE approach	Bayesian approach
Assume that there is an unknown but fixed parameter θ .	Represents uncertainty about the unknown parameter θ .
Estimate θ with some confidence.	Uses probability to quantify this uncertainty. Unknown parameters as random variables.
Prediction using the estimated parameter value.	Prediction follows from the rules of probability- expectation over the unknown parameters.

Table 8.8: The differences between the MLE and Bayesian approaches

Dirichlet's function:

$$P(\theta) = \frac{(\sum_{j=1}^k \alpha_j - 1)!}{(\alpha_1 - 1)! (\alpha_2 - 1)! \dots (\alpha_k - 1)!} \cdot \prod_{j=1}^k \theta_j^{\alpha_j - 1}$$

This is the density function of a k-dimensional random variable θ with hyperparameters $\alpha_1, \dots, \alpha_k$. If we know these hyperparameters, we will know the density function explicitly as a function of $\theta_1, \dots, \theta_k$. Note that we can "drop" the coefficient (it is a normalization constant) and get:

$$P(\theta) \propto \prod_{j=1}^k \theta_j^{\alpha_j - 1}$$

We now use the Dirichlet distribution as a prior. Recall that the likelihood function for multinomial data is:

$$L(\theta : D) = P(D | \theta) = \prod_{j=1}^k \theta_j^{N_j}$$

From equation 8.12 we get:

$$P(\theta | D) \stackrel{(*)}{\propto} P(\theta) P(D | \theta) \propto \prod_{j=1}^k \theta_j^{\alpha_j - 1} \prod_{j=1}^k \theta_j^{N_j} = \prod_{j=1}^k \theta_j^{\alpha_j + N_j - 1} \quad (8.13)$$

We can see that the posterior density function $P(\theta | D)$ has the same form as the prior $P(\theta)$, with hyperparameters $\alpha_j + N_j - 1, \dots, \alpha_k + N_k - 1$. Therefore we can conclude that Dirichlet is **closed** under multinomial sampling (i.e. both the prior and posterior distributions are Dirichlet). Note that in the transition marked (*) we disregarded the denominator $P(D)$ from equation 8.12. This denominator and the normalization constant, which was "dropped" from

Dirichlet's function, give us the normalization constant of our new Dirichlet distribution. Since the posterior is Dirichlet, we get:

$$P(x[m+1] = j \mid D) \stackrel{(*)}{=} \int \theta_j \cdot P(\theta \mid D) d\theta \stackrel{(**)}{=} \frac{\alpha_j + N_j}{\sum_l (\alpha_l + N_l)}$$

where $(*)$ is from equation 8.11 and $(**)$ is the result of solving the integral with $P(\theta \mid D)$, which appears in equation 8.13. We can learn from the above expressions that α indicates the "degree of influence" the prior has on the posterior result. The larger α is, the greater the effect the prior has, because it means we have simulated a situation with more experiments (for each j we have $\alpha_j + N_j$ experiments instead of the N_j we had in the likelihood function) and our "added" experiments α_j have a greater part out of the whole $\alpha_j + N_j$ experiments. We can now generalize:

$$P(x_i[m+1] = j \mid pa_i, D) = \frac{\alpha(x_i = j, pa_i) + N(x_i = j, pa_i)}{\sum_l (\alpha(x_i = l, pa_i) + N(x_i = l, pa_i))} = \frac{\alpha(x_i = j, pa_i) + N(x_i = j, pa_i)}{\alpha(pa_i) + N(pa_i)}$$

Learning Multinomial Parameters: Summary

We count $N(x_i, pa_i)$ and get the following estimations:

$$\underbrace{\hat{\theta}_{x_i|pa_i} = \frac{N(x_i, pa_i)}{N(pa_i)}}_{\text{MLE}} \quad \underbrace{\hat{\theta}_{x_i|pa_i} = \frac{\alpha(x_i, pa_i) + N(x_i, pa_i)}{\alpha(pa_i) + N(pa_i)}}_{\text{Bayesian (Dirichlet)}}$$

Both can be implemented in an on-line manner by accumulating counts.

8.4.11 Learning Structure

As described in section 8.4.2 (The Learning Problem), given a training set D , our goal is to find a network that best matches D . This time, however, the structure is unknown and we wish to learn it on top of learning the parameters. We therefore describe the following **optimization problem**:

Given an input of:

- Training data
- A scoring function (including priors) to measure our suggested network's match with the data
- Set of possible structures- a set of networks (DAGs, trees, etc) from which we need to find the best structure

we need to output the network (or networks) that maximize the score.

As before, we would like to take advantage of the decomposability property, and therefore we look for a scoring function such that the score of the network will be the sum of the nodes' scores. An example for such a score is the **BDE score**:

$$P(G \mid D) \propto P(D \mid G)P(G) \stackrel{(*)}{=} P(G) \int P(D \mid G, \theta)P(\theta \mid G)d\theta$$

where (*) holds due to the total probability theorem for $P(D \mid G)$. Following particular assumptions, the score satisfies the decomposability property:

$$Score(G : D) = \sum_i Score(x_i \mid Pa_i^G : D)$$

The problem of finding the optimal structure is NP-hard and we therefore address the problem by using heuristic searching. We traverse the space of possible networks, looking for high-scoring structures. This can be done by using searching techniques, such as:

- Simulated Annealing - we start with a certain structure and perform random modifications to it. The number or size of modifications is reduced with time so we start with a search of the entire solution space and later use smaller modifications in order to achieve a local minima.
- Greedy hill-climbing- we start with a certain structure and we add modifications until we reach a local maximum. In this technique, we save time by not going over all the possibilities.

The typical operations performed in our heuristic search can be seen in figure 8.13- we can add, remove or reverse the direction of an edge, thus creating a new structure, which has a new score we can calculate. Traversing all possible edges and applying the above modifications, we can find the highest-scoring new structure and reiterate the process from it. The number of edge modifications in each step of the greedy hill-climbing algorithm is $O(n^2)$, where n is the number of nodes (because there are $O(n^2)$ edges). Once again, the importance of the decomposability property is demonstrated- each modification of a certain edge is local to the node that particular edge enters. The decomposability enables us to update the score of either only one node (in case of an addition or removal of an edge) or two nodes (in case of a reverse in an edge's direction), without having to recalculate the entire network's score.

An algorithm for learning Bayesian network's structure

An example for a structure inference algorithm is the algorithm by Pen and Ding [21]. Their algorithm uses local learning in order to generate parts of the Bayesian network and merge

them into a valid network. This method is based on the locality of the network's score 11.14. The algorithm is composed of 3 steps :

- Generation of local structure. For each node, the optimal set of parents is found by starting with an optimal (i.e. likelihood maximizing) singleton set and adding nodes in order to improve the solution (i.e. increase the likelihood for that node)
- Elimination of cycles. After composing the local structures, found in the first step, into a Bayesian network the algorithm makes sure the network is valid. The network might contain cycles which have to be removed in order to create a valid network (Bayesian networks have to be DAGs¹). The challenge is to remove the cycles with minimal damage to the network's score.
- Local structure perturbation. To asses and improve the stability of the network the algorithm creates local perturbations (e.g. eliminating a random edge and replacing it by other local edges) and tests for improvement of the network's score. Unstable edges whose replacement improves the total score are replaced.

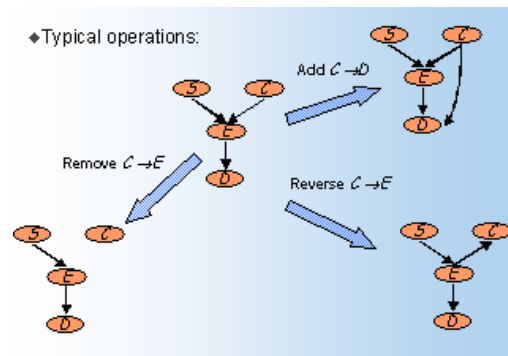


Figure 8.13: The typical operations performed in the heuristic search

8.5 Applying Bayesian Networks to Expression Data

This section describes an approach for analyzing gene expression data using Bayesian network learning techniques. We model the expression level of each gene as a random variable. In addition, other attributes that affect the system can be modeled as random variables. These can include a variety of attributes of the sample, such as experimental conditions, temporal indicators (i.e., the time/stage that the sample was taken from), background variables (e.g., which clinical procedure was used to get a biopsy sample), and exogenous cellular conditions.

¹Direct Acyclic Graph

By learning a Bayesian network based on the statistical dependencies between these variables, we can answer a wide range of queries about the system. For example, does the expression level of a particular gene depend on the experimental condition? Is this dependence direct, or indirect? If it is indirect, which genes mediate the dependency?

We now describe how one can learn such a model from expression data. Many important issues arise when learning in this domain. These involve statistical aspects of interpreting the results, algorithmic complexity issues in learning from the data, and preprocessing the data.

Most of the difficulties in learning from expression data revolve around the following central point: Contrary to previous applications of learning Bayesian networks, expression data involves transcript levels of thousands of genes while current datasets contain at most a few dozen samples. This raises problems in computational complexity and the statistical significance of the resulting networks. On the positive side, genetic regulation networks are sparse, i.e., given a gene, it is assumed that no more than a few dozen genes directly affect its transcription. Bayesian networks are especially suited for learning in such sparse domains.

8.5.1 Network Features

When learning models with many variables, small datasets are not sufficiently informative to significantly determine that a single model is the "right" one. Instead, many different networks should be considered as reasonable explanation of the given the data. From a Bayesian perspective, we say that the posterior probability over models is not dominated by a single model (or equivalence class of models). We would like to analyze this set of plausible (i.e., high-scoring) networks. Although this set can be very large, we might attempt to characterize features that are common to most of these networks, and focus on learning them.

Before we examine the issue of inferring such features, we briefly discuss two classes of features involving pairs of variables:

Markov Relations

A relation of this type specifies if Y is in the Markov blanket of X , where the Markov blanket of X is the minimal set of variables that shield X from the rest of the variables in the model (see Figure 8.14 for an example). More precisely, X given its Markov blanket is independent from the remaining variables in the network. It is easy to check that this relation is symmetric: Y is in X 's Markov blanket if and only if there is either an edge between them, or both are parents of another variable [19]. In the context of gene expression

analysis, a Markov relation indicates that the two genes are related in some joint biological interaction or process. Note, that two variables in a Markov relation are directly linked in the sense that no variable in the model mediates the dependence between them. It remains possible that an unobserved variable (e.g., protein activation) is an intermediate in their interaction.

Order Relations

An order relation specifies if X is an ancestor of Y in all the networks of a given equivalence class. That is, if the given PDAG¹ contains a directed path from X to Y . This type of relation does not involve only a close neighborhood, but rather captures a global property. Thus, we view such a relation as an indication, rather than evidence, that X might be a causal ancestor of Y .

While at this point we handle only pairwise features, it is clear that this analysis is not restricted to them, and we should examine also features that are more complex (see [20]).

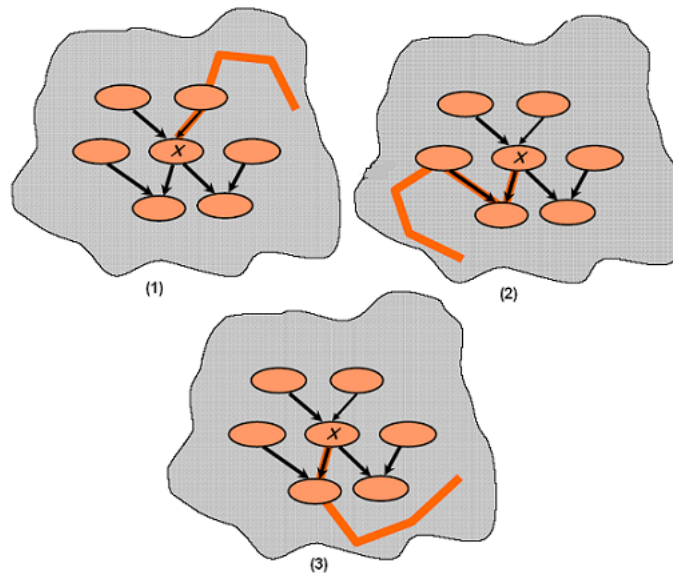


Figure 8.14: Source [26]. An example for a Markov blanket. This Markov blanket of X contains all paths from X to other nodes. There are three kinds of such paths as shown in the figure: (1) *Upward paths* the parents of X . (2) *Sideway paths* the spouses of X . (3) *Downward paths* the children of X .

¹Partially Directed Acyclic Graph - A Graph in which only a part of the edges are directed

8.5.2 Estimating Statistical Confidence in Features

We now face the following problem: To what extent do the data support a given feature? More precisely, we want to estimate a measure of confidence in the features of the learned networks, where "confidence" approximates the likelihood that a given feature is actually true (i.e., is based on a genuine correlation and causation) (see Figure 8.16).

An effective, and relatively simple, approach for estimating confidence is the bootstrap method [6]. The main idea behind the bootstrap is simple. We generate "perturbed" versions of the original dataset, and learn from them. In this way we collect many networks, all of which are fairly reasonable models of the data. These networks show how small perturbations to the data can affect many of the features. In our context, we use the bootstrap as follows:

- For $i = 1 \dots m$ (in the experiments, we set $m = 200$):
 - Resample with replacement N instances from D . Denote by D_i the resulting dataset.
 - Apply the learning procedure on D_i to deduce a network structure \hat{G}_i .
- For each feature f of interest calculate

$$conf(f) = \frac{1}{m} \sum_{i=1}^m f(\hat{G}_i) \quad (8.14)$$

where $f(G)$ is 1 if f is a feature in G , and 0 otherwise.

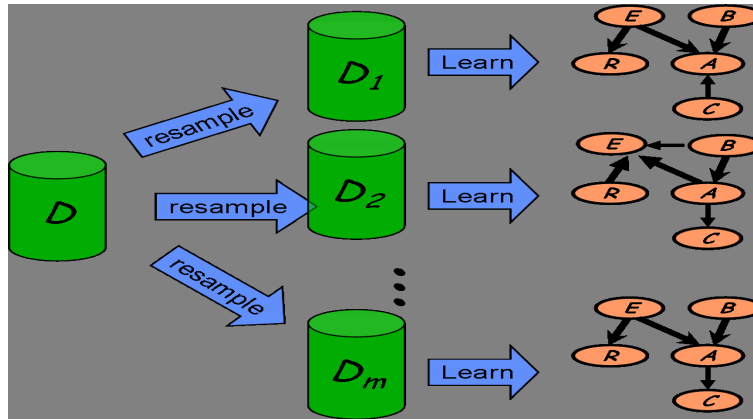


Figure 8.15: A Bootstrap method. Only features with high confidence will be accepted

8.5.3 Efficient Learning Algorithms

In section 11.4 we formulated learning Bayesian network structure as an optimization problem in the space of directed acyclic graphs. The number of such graphs is super-exponential in the number of variables. As we consider hundreds and thousands of variables, we must deal with an extremely large search space. Therefore, we need to use (and develop) efficient search algorithms.

To facilitate efficient learning, we need to be able to focus the attention of the search procedure on relevant regions of the search space, giving rise to the sparse candidate algorithm [12]. The main idea of this technique is that we can identify a relatively small number of candidate parents for each gene based on simple local statistics (such as correlation). We then restrict our search to networks in which only the candidate parents of a variable can be its parents, resulting in a much smaller search space in which we can hope to find a good structure quickly.

8.6 Experimental Results

The Bayesian Networks approach was applied by Freidman et al. [9] to two datasets: the data of Spellman et al. [23] and the data of Hughes et al. [15]. From this point and on, we will refer only to the data of Spellman et al. We refer the reader to [20] for details about the results from the data of Hughes et al.

The data contains 79 gene expression measurements of the mRNA levels of 6177 *S. cerevisiae* ORFs. These experiments measure expression in fixed time intervals under different cell cycle synchronization methods. Spellman et al. identified 800 genes whose expression varied over the different cell-cycle stages. They clustered these 800 genes, based on the similarity of expression profiles, resulting 8 major clusters, which contained 250 genes in total. The variables of the learned networks were the expression level of each of these 800 genes. Some of the robustness analysis was performed only on the set of 250 genes that appear in the 8 major clusters.

Freidman et al. [9] used the Sparse Candidate algorithm with a 200-fold bootstrap in the learning process. The learned features show that intricate structure can be recovered even from such small data sets. It is important to note that this learning algorithm uses no prior biological knowledge nor constraints. All learned networks and relations are based solely on the information conveyed in the measurements themselves. These results are available in [11]. Figure 8.16 illustrates the graphical display of results of this analysis.

8.6.1 Robustness Analysis

Freidman et al. performed a number of tests to analyze the statistical significance and robustness of their procedure. They carried most of these tests on the smaller 250 gene data set

for computational reasons. To test the credibility of their confidence assessment, they created a random data set by randomly permuting the order of the experiments independently for each gene. Thus for each gene the order was random, but the composition of the series remained unchanged. In such a data set, genes are independent of each other, and thus we do not expect to find "real" features. As expected, both order and Markov relations in the random data set have significantly lower confidence. Clearly, the distribution of confidence estimates in the original data set have a longer and heavier tail in the high confidence region. The runs on the random data sets do not learn almost anything with a confidence level above 0.8, which can lead us to believe that most features that are learned in the original data set with such confidence levels originate in true signals in the data. Also, the confidence distribution for the real dataset is concentrated closer to zero than the random distribution. This suggests that the networks learned from the real data are sparser. Since the analysis was not performed on the whole *S. cerevisiae* genome, Freidman et al. also tested the robustness of their analysis to the addition of more genes, comparing the confidence of the learned features between the 250 and 800 gene datasets.

8.6.2 Biological Analysis

Freidman et al. believe that the results of this analysis can be indicative of biological phenomena in the data. This is confirmed by their ability to predict sensible relations between genes of known function. We now examine several consequences from this analysis. We consider, in turn, the order relations and Markov relations found.

Order Relations

The most striking feature of the high confidence order relations, is the existence of dominant genes. Out of all 800 genes only few seem to dominate the order (i.e., appear before many genes). The intuition is that these genes are indicative of potential causal sources of the cell-cycle process. A list of the highest scoring dominating genes appears in Table 8.9.

Inspection of the list of dominant genes reveals quite a few interesting features. Among the dominant genes are those directly involved in cell-cycle control and initiation. For example, CLN1, CLN2 and CDC5, whose functional relation has been established [4, 5]. Other genes, like MCD1 and RFA2, were found to be essential [14]. These are clearly key genes in basic cell functions, involved in chromosome dynamics and stability (MCD1) and in nucleotide excision repair (RFA2). Most of the dominant genes encode nuclear proteins, and some of the unknown genes are also potentially nuclear: (e.g., YLR183C contains a forkhead associated domain which is found almost entirely among nuclear proteins). Some of them are components of pre-replication complexes. Others (like RFA2, POL30 and MSH6) are

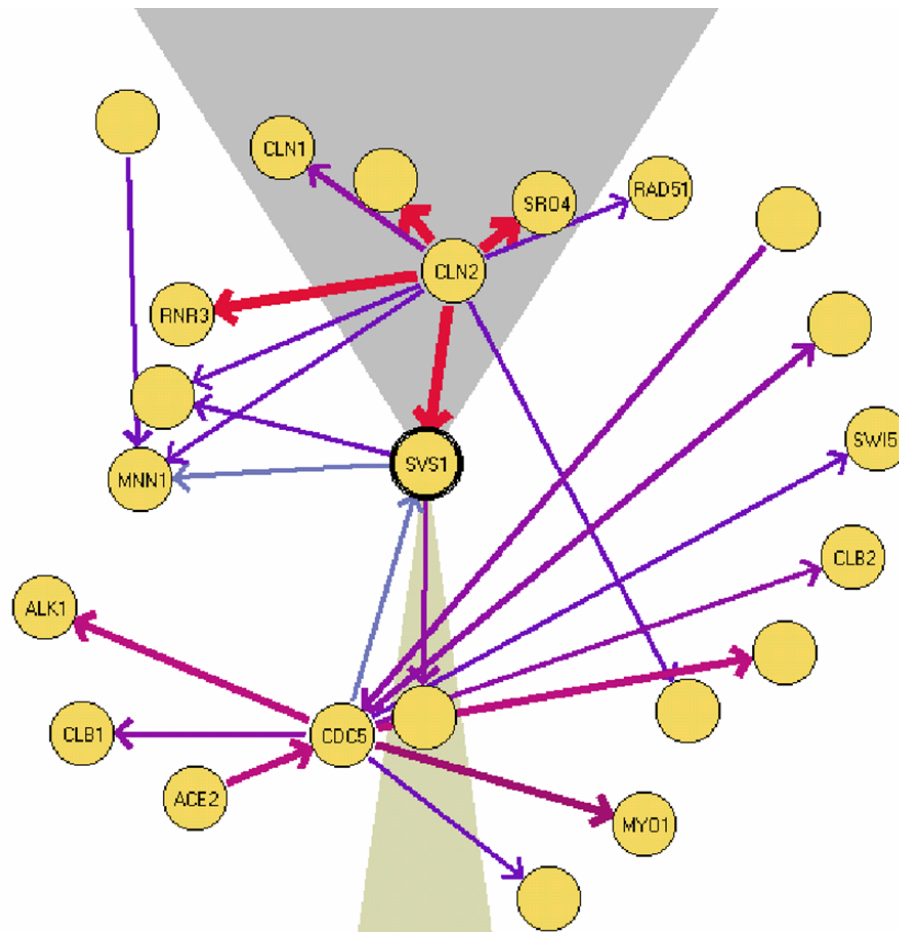


Figure 8.16: Source: [9]. An example of the graphical display of Markov features. This graph shows a "local map" for the gene SVS1. The width (and color) of edges corresponds to the computed confidence level. An edge is directed if there is a sufficiently high confidence in the order between the pair genes connected by the edge. This local map shows that CLN2 separates SVS1 from several other genes. Although there is a strong connection between CLN2 to all these genes, there are no other edges connecting them. This indicates that, with high confidence, these genes are conditionally independent given the expression level of CLN2.

Gene/ORF	Dominance Score	# of desc. genes		Notes
		> .8	> .7	
YLR183C	551	609	708	Contains forkheaded associated domain, thus possibly nuclear
MCD1	550	599	710	Mitotic chromosome determinant, null mutant is inviable
CLN2	497	495	654	Role in cell cycle START, null mutant exhibits G1 arrest
SRO4	463	405	639	Involved in cellular polarization during budding
RFA2	456	429	617	Involved in nucleotide excision repair, null mutant is inviable
YOL007C	444	367	624	Glycophospholipid surface protein, Null mutant is slow growing
GAS1	433	382	586	
YOX1	400	243	556	Homeodomain protein that binds leutRNA gene
YLR013W	398	309	531	Required for DNA replication and repair, Null mutant is inviable
POL30	376	173	520	
RSR1	352	140	461	GTPbinding protein of the ras family involved in bud site selection
CLN1	324	74	404	Role in cell cycle START, null mutant exhibits G1 arrest
YBR089W	298	29	333	Required for mismatch repair in mitosis and meiosis
MSH6	284	7	325	

Table 8.9: Source [9]. List of dominant genes in the ordering relations (top 14 out of 30). The first column specifies the name of the gene/ORF, the second column specifies the level of dominance score of the gene/ORF as appeared in the experiments results, the next column contains the number of descendent genes with a level of confidence higher than 0.8, the next column contains the number of descendent genes with a level of confidence higher than 0.7 and the last column supplies additional biological information about the gene/ORF.

involved in DNA repair. It is known that DNA repair is a prerequisite for transcription, and DNA areas which are more active in transcription, are also repaired more frequently [18, 25]. A few non nuclear dominant genes are localized in the cytoplasm membrane (SRO4 and RSR1). These are involved in the budding and sporulation process which have an important role in the cell-cycle. RSR1 belongs to the ras family of proteins, which are known as initiators of signal transduction cascades in the cell.

Markov Relations

Inspection of the top Markov relations reveals that most are functionally related. A list of the top scoring relations can be found in Table 8.10. Among these, all involving two known genes make sense biologically. When one of the ORFs is unknown careful searches using Psi-Blast [1], Pfam [22] and Protomap [27] can reveal firm homologies to proteins functionally related to the other gene in the pair (e.g. YHR143W, which is paired to the endochitinase CTS1, is related to EGT2 - a cell wall maintenance protein). Several of the unknown pairs are physically adjacent on the chromosome and, thus, are presumably regulated by the same mechanism (see [2]), although special care should be taken for pairs whose chromosomal location overlap on complementary strands, since in these cases we might see an artifact resulting from cross-hybridization. There are some interesting Markov relations found that are not discovered using clustering techniques. One such regulatory link is FAR1-ASH1: both proteins are known to participate in a mating type switch. The correlation of their expression patterns is low and [23] cluster them into different clusters. Among the high confidence markov relations, one can also find examples of conditional independence, i.e., a group of highly correlated genes whose correlation can be explained within the resulted network structure. One such example involves the genes: CLN2, RNR3, SVS1, SRO4 and RAD41. Their expression is correlated and in [23] all appear in the same cluster. In the resulting network CLN2 is with high confidence a parent of each of the other 4 genes, while no links are found between them. This suits biological knowledge: CLN2 is a central and early cell cycle control, while there is no clear biological relationship between the others.

Confidence	Gene 1	Gene 2	Notes
1.0	YKL163WPIR3	YKL164CPIR1	Close locality on chromosome
0.985	PRY2	YKR012C	No homolog found
0.985	MCD1	MSH6	Both bind to DNA during mitosis
0.98	PHO11	PHO12	Both nearly identical acid phosphatases
0.975	HHT1	HTB1	Both are Histones
0.97	HTB2	HTA1	Both are Histones
0.94	YNL057W	YNL058C	Close locality on chromosome
0.94	YHR143W	CTS1	Homolog to EGT2 cell wall control, both do cytokinesis
0.92	YOR263C	YOR264W	Close locality on chromosome
0.91	YGR086	SIC1	Homolog to mammalian nuclear ran protein, both involved in nuclear function
0.9	FAR1	ASH1	Both part of a mating type switch, expression uncorelated
0.89	CLN2	SVS1	Function of SVS1 unknown, possible regulation mediated through
SWI6	0.88	YDR033W	NCE2 Homolog to transmembrane proteins, suggesting both involved in protein secretion
0.86	STE2	MFA2	A mating factor and receptor
0.85	HHF1	HHF2	Both are Histones
0.85	MET10	ECM17	Both are sulfite reductases
0.85	CDC9	RAD27	Both participate in Okazaki fragment processing

Table 8.10: Source [9]. List of top Markov relations. The first column describes the level of confidence of the relation, the next two columns contain the names of the two genes of the relation and the last column supplies additional biological information referring the relation.

8.7 Using biological knowledge to increase robustness

An interesting method to increase the robustness of created Bayesian networks is based on the usage of prior biological knowledge regarding the structure of the Bayesian network. This knowledge is used to filter possible networks created in the Bootstrap phase. in [24] the Bayesian network is limited to a two level tree (see 8.17).

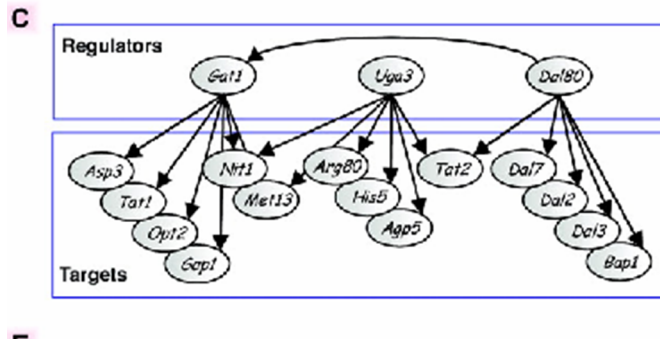


Figure 8.17: In [24] the structure of the Bayesian network is limited to a two level tree, based on the naive separation between regulators and targets

In [16] a more sophisticated structure is imposed, in which the network is composed of modules - sets of co-regulated parameters (i.e. immediate descendants of the same parameters) (see 8.18).

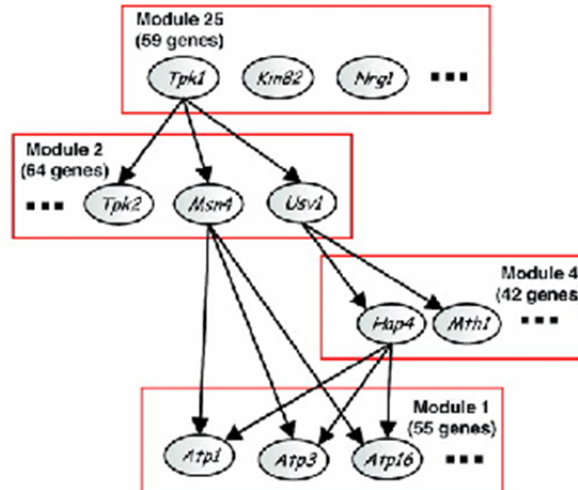


Figure 8.18: In [16] the structure of the Bayesian network is limited to a tree of modules

8.8 Improving The Framework

The framework we described here can be expanded in a number of promising directions:

- Correct handling of hidden variables (active proteins)
- Developing the theory for learning local probability models that are capable of dealing with the continuous nature of the data.
- Improving the theory and algorithms for estimating confidence levels.
- Incorporating biological knowledge (such as possible regulatory regions) as prior knowledge to the analysis.
- Improving the search heuristics.
- Applying Dynamic Bayesian Networks ([10]) to temporal expression data.
- Dealing with feedback loops which are common to metabolic pathways

Finally, one of the most exciting longer term prospects of this line of research is discovering causal patterns from gene expression data. We can build on and extend the theory for learning causal relations from data and apply it to gene expression. The theory of causal networks allows learning both from observational data and interventional data, where the experiment intervenes with some causal mechanisms of the observed system. In gene expression context, we can model knockout/overexpressed mutants as such interventions. Thus, we can design methods that deal with mixed forms of data in a principled manner (See [3] for a recent work in this direction). In addition, this theory can provide tools for experimental design, that is, understanding which interventions are deemed most informative to determining the causal structure in the underlying system. Friedman et al. have extended their framework in this direction (see [20]).

Bibliography

- [1] S.F. Altschul, T.L. Madden, A.A. Schaffer, J. Zhang, Z. Zhang, W. Miller, and D.J. Lipman. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic Acids Res.*, 25, 1998.
- [2] T. Blumenthal. Gene clusters and polycistronic transcription in eukaryotes. *Bioessays*, 480-487, 1998.
- [3] G. Cooper and C. Yoo. Causal discovery from a mixture of experimental and observational data. *Proc. Fifteenth Conference on Uncertainty in Artificial Intelligence (UAI'99)*, pages 116–125, 1999.
- [4] F. Cvrckova and K. Nasmyth. Yeast g1 cyclins cln1 and cln2 and a gap-like protein have a role in bud formation. *EMBO*, pages 12:5277–5286, 1993.
- [5] M.A. Drebot, G.C. Johnston, J.D. Friesen, and R.A. Singer. An impaired rna polymerase ii activity in *saccharomyces cerevisiae* causes cell-cycle inhibition at start. *Mol Gen Genet*, pages 241:327–334, 1993.
- [6] B. Efron and R.J. Tibshirani. *An Introduction to the Bootstrap*. Chapman and Hall, London, 1993.
- [7] J. Finn. *An Introduction to Bayesian Networks*. UCL Press, 1996.
- [8] N. Friedman and D. Koller. <http://www.cs.huji.ac.il/~nir/nips01-tutorial/>.
- [9] N. Friedman, M. Linial, I. Nachman, and D. Pe'er. Using bayesian networks to analyze expression data. *Journal of Computational Biology*, pages 7:601–620, 2000.
- [10] N. Friedman, K. Murphy, and S. Russell. Learning the structure of dynamic probabilistic networks. *Proc. Fourteenth Conference on Uncertainty in Artificial Intelligence (UAI'98)*, pages 139–147, 1998.
- [11] N. Friedman, I. Nachman, and D. Pe'er. <http://www.cs.huji.ac.il/labs/compbio/expression/>.

- [12] N. Friedman, I. Nachman, and D. Pe'er. Bayesian network structure from massive datasets: The 'sparse candidate' algorithm. *Proc. Fifteenth Conference on Uncertainty in Artificial Intelligence (UAI'99)*, pages 196–205, 1999.
- [13] C.M. Grinstead and J.L. Snell. *Introduction to Probability*. American Mathematical Society, second edition, 1997.
- [14] V. Guacci, D. Koshland, and A. Strunnikov. A direct link between sister chromatid cohesion and chromosome condensation revealed through the analysis of *mcd1* in *s. cerevisiae*. *Cell*, pages 91:47–57, 1997.
- [15] T.R. Hughes et al. Functional discovery via a compendium of expression profiles. *Cell*, pages 102:109–126, 2000.
- [16] D. Koller, S.K. Kim, J.M. Stuart, and E. Segal. A gene-coexpression network for global discovery of conserved genetic modules. *Science*, pages 302:249–255, 2003.
- [17] T. Leviatan and A. Raviv. *Introduction to Probability and Statistics- Statistical Inference*. Amichai Publishing House, 1997.
- [18] W. G. McGregor. Dna repair, dna replication, and uv mutagenesis. *Investig Dermatol Symp Proc*, pages 4:1–5, 1999.
- [19] J. Pearl. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, San Francisco, 1988.
- [20] D. Pe'er, A. Regev, G. Elidan, and N. Friedman. Inferring subnetworks from perturbed expression profiles. *Proc. Ninth International Conference on Intelligent Systems for Molecular Biology (ISMB '01)*, 2001.
- [21] H. Peng and C. Ding. Structure search and stability enhancement of bayesian networks. *Proceeding of the Third IEEE Conference on Data Mining*, 2003.
- [22] E.L. Sonnhammer, S.R. Eddy, E. Birney, A. Bateman, and R. Durbin. Pfam: multiple sequence alignments and hmm-profiles of protein domains (<http://pfam.wustl.edu/>). *Nucleic Acids Res*, pages 26:320–322, 1998.
- [23] P.T. Spellman et al. Comprehensive identification of cell cycle-regulated genes of the yeast *saccharomyces cerevisiae* by microarray hybridization. *Molecular Biology of the cell*, pages 9:3273–3297, 1998.
- [24] A. Tanay, D. Pe'er, and A. Regev. A fast and robust method to infer and characterize an active regulator set for molecular pathways. *Bioinformatics*, 18(1).

- [25] S. Tornaletti and P. C. Hanawalt. Effect of dna lesions on transcription elongation. *Biochimie*, pages 81:139–146, 1999.
- [26] Y. Weiss. <http://www.cs.huji.ac.il/~pmai/tirguls/tirgul3.pdf>.
- [27] G. Yona, N. Linial, and M. Linial. Protomap - automated classification of all protein sequences: a hierarchy of protein families, and local maps of the protein space. *Proteins: Structure, Function, and Genetics*, pages 37:360–378, 1998.