

Faculty of Computer Science, Dalhousie University

18-Nov-2025

CSCI 4152/6509 — Natural Language Processing**Lecture 14: P0 Discussion (2); Bayesian Networks and HMM Inference**

Location: Studley LSC-Psychology P5260

Instructor: Vlado Keselj

Time: 14:35 – 15:55

Previous Lecture

- P0 discussion (1): P-03, P-07, P-08, P-09, P-11, P-12, P-20, P-21, P-22, P-23
- HMM Brute-force approach
- HMM Inference: Viterbi algorithm

P0 Topics Discussion (2)

- Continued discussion of individual projects as proposed in P0 submissions (part 2)
- Projects discussed: P-01, P-04, P-06, P-13, P-14, P-15, P-16, P-17

16 HMM as Bayesian Network*Slide notes:***HMM as Bayesian Network**

- **Viterbi** algorithm is an **efficient** way to solve a **special** problem:
 - completion with known observables and unknown hidden nodes of an HMM
- **General** approach:
 - Treat HMM as **Bayesian Network**
 - Apply **Product-Sum** (i.e., “Message-passing”) algorithm for efficient inference

In the previous section, we saw how we can do efficient inference in an HMM model using the Viterbi algorithm. The inference task that we were solving was a specialized completion task where we knew values of all observable variables, and we calculated most likely values of all hidden variables. In order to solve more general inference tasks, such as marginalization, or for example completion with arbitrary known variables, we would need to modify the Viterbi algorithm in a non-trivial way, or possibly come up with a significantly different new algorithm.

A much more general and flexible way to approach inference in HMMs is to treat them as Bayesian Networks, since HMM is a specialized version of a more general model called *Bayesian Network*.

16.1 Bayesian Networks

Note regarding course scope: Some detailed discussion in this section is not completely covered in the lectures, so that part will not be used in assignments or the exam. For the purposes of the course, it is sufficient to understand how product-sum algorithms work and be able to apply them or implement them.

Slide notes:

Bayesian Network Model

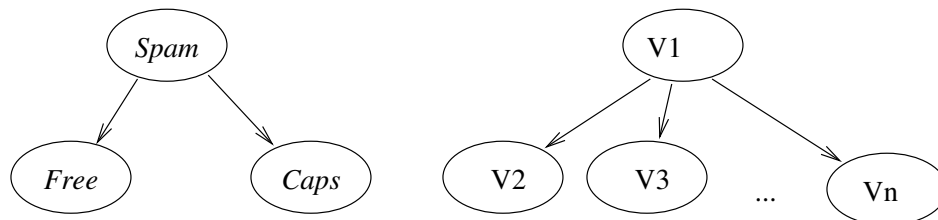
- Also known as: Belief Networks, or Bayesian Belief Networks
- A directed acyclic graph (DAG)
 - Each node representing a random variable
 - Edges representing causality (probabilistic meaning)
- Conditional Probability Table (CPT) for each node
- Bayesian Network assumption:

$$P(\underbrace{V_1, V_2, \dots, V_n}_{\text{full configuration}}) = \prod_{i=1}^n P(V_i | \mathbf{V}_{\pi(i)})$$

The Viterbi algorithm provides an efficient way to solve a particular inference problem in the HMM model. It is the completion problem of finding the most likely hidden states given the observable values. If we represent the HMM model in the unrolled graphical form, we can see that this inference problem is only a special case of the inference problem, and that we may want to calculate other inference formulas, including different marginalization, conditioning, or completion cases. Many of these problems can be efficiently solved if we see HMM as a special case of a more general concept of a Bayesian Network.

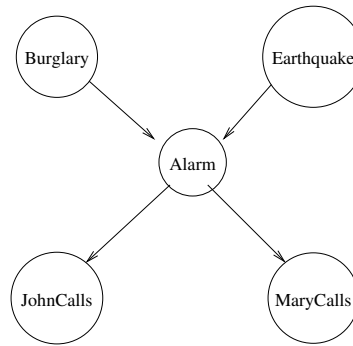
Bayesian Networks (also known as **belief networks**, **Bayesian belief networks**, or **decision networks**) provide a way to create structured probabilistic models in a more general way. We saw previously that between the joint distribution model and the fully independent model we have useful models such as the Naïve Bayes model, the Ngram or Markov chain model, and the Hidden Markov Model. All of these models had a graphical representation, and the Bayesian Network model makes this more general by providing a model associated with any directed acyclic graph.

The following graphs shows dependence structure for the Naïve Bayes model of the spam detection, and for the Naïve Bayes model in general:



Example

To introduce the general concept of Bayesian Networks, we will consider the known Burglar-Earthquake example, frequently used in the literature (e.g., the AI textbook by Russell and Norvig):



We can follow a topological sort of the above graph and use the chain rule for the conditional probability, which can be derived from the definition of the conditional probability, to obtain the equation:

$$P(B, E, A, J, M) = P(B)P(E|B)P(A|B, E)P(J|A, B, E)P(M|J, A, B, E)$$

The graph denotes some dependence assumptions, which include: B and E are independent variables, hence $P(E|B) = P(E)$, J depends only on A, hence $P(J|A, B, E) = P(J|A)$, and M depends only on A, hence $P(M|J, A, B, E) = P(M|A)$. If we make these substitutions in the above equation, we obtain the following, **Bayesian network assumption**:

$$P(B, E, A, J, M) = P(B)P(E)P(A|B, E)P(J|A)P(M|A)$$

This assumption implies that to evaluate any probability of a complete configuration of the model we need to keep only the parameters for the following conditional probabilities: $P(B)$, $P(E)$, $P(A|B, E)$, $P(J|A)$, and $P(M|A)$. These parameters are given in corresponding **conditional probability tables (CPTs)**.

Let us assume that the conditional tables for our example are:

				B	E	A	$P(A B, E)$									
				T	T	T	0.95									
				T	T	F	0.05				A	J	$P(J A)$	A	M	$P(M A)$
B	$P(B)$	E	$P(E)$	T	F	T	0.94	T	T	0.90	T	T	0.70			
T	0.001	T	0.002	T	F	F	0.06	T	F	0.10	T	F	0.30			
F	0.999	F	0.998	F	T	T	0.29	F	T	0.05	F	T	0.01			
				F	T	F	0.71	F	F	0.95	F	F	0.99			
				F	F	T	0.001									
				F	F	F	0.999									

The following is a formal definition of Bayesian Network:

Definition 16.1 (Bayesian Network) A Bayesian Network is defined by a directed acyclic graph (DAG) and a collection of conditional probability tables, where nodes in the graph represent random variables and directed edges in the graph represent conditional independence assumptions. The edges represent causality among the random variable nodes so that there is an edge from a variable to another variable if and only if there is a direct dependence between these variables. If V_j ($1 \leq j \leq n$) are the random variables and $\mathbf{V}_{\pi(j)}$ are parent variables of any variable V_j , i.e., then the Bayesian Network assumption can be written in the following way:

$$P(\text{all variables}) = \prod_{j=1}^n P(V_j | \mathbf{V}_{\pi(j)})$$

Hence, a Bayesian network consists of two components: a directed graph and a set of lookup tables for conditional probabilities for each variable V_i . If $\pi(i)$ are parent nodes of V_i , then for all possible values of variables $\mathbf{V}_{\pi(i)}$ and V_i , the conditional probability table (CPT) specifies the probability of $P(V_i|\mathbf{V}_{\pi(i)})$.

Example (Burglar-Earthquake cont'd)

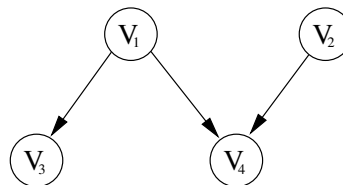
The tables for the Burglar-Earthquake example:

B	$P(B)$	E	$P(E)$	B	E	A	$P(A B, E)$	A	J	$P(J A)$	A	M	$P(M A)$
T	0.001	T	0.002	T	T	T	0.95	T	T	0.90	T	T	0.70
F	0.999	F	0.998	T	T	F	0.05	T	F	0.10	T	F	0.30
				T	F	T	0.94	F	T	0.05	F	T	0.01
				T	F	F	0.06	F	F	0.95	F	F	0.99
				F	T	T	0.29						
				F	T	F	0.71						
				F	F	T	0.001						
				F	F	F	0.999						

Number of Parameters

If the number of parent nodes is k , and if we assume that each variable may have m different values, then the conditional probability table has m^{k+1} rows. For each of m^k combinations of values of parent nodes, there is one constraint on the probability distribution $\sum_x P(V_i = x|\mathbf{V}_{\pi(i)}) = 1$, which will result in m^k constraints; i.e., there are $m^{k+1} - m^k = (m - 1)m^k$ free parameters. If the maximal number of parents for a whole network is k , then the total number of free parameters is not greater than $n(m - 1)m^k$.

Example. For example, let us calculate the number of free parameters of the following Bayesian Network:



The Bayesian assumption for the network above is:

$$\begin{aligned}
 &P(V_1 = x_1, V_2 = x_2, V_3 = x_3, V_4 = x_4) \\
 &= P(V_1 = x_1) P(V_2 = x_2) P(V_3 = x_3|V_1 = x_1) P(V_4 = x_4|V_1 = x_1, V_2 = x_2)
 \end{aligned}$$

How many parameters are needed to represent the network?

For each variable store a conditional probability table of size

$$m \cdot m^{\# \text{parents}} - m^{\# \text{parents}} \quad (\text{constraints})$$

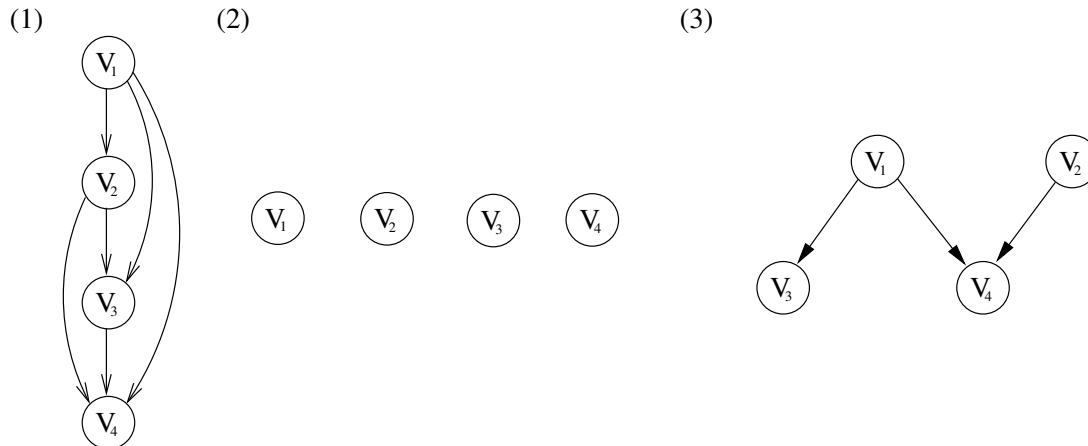
For the above network, the number of free parameters is:

$$\begin{aligned}
 &m + m + m^2 + m^3 \quad \text{parameters} \\
 &- 1 - 1 - m - m^2 \quad \text{constraints} \\
 &= m^3 + m - 2
 \end{aligned}$$

Representational Power

Using the Bayesian network model we can represent many other models: the full joint distribution model, fully independent model, Naïve Bayes model, Hidden Markov Model, and other structured models.

Some examples of the Bayesian Networks are:



Number of parameters in each model:

$$(1) \quad (m-1) + (m^2 - m) + (m^3 - m^2) + (m^4 - m^3) = m^4 - 1$$

$$(2) \quad (m-1) + (m-1) + (m-1) + (m-1) = 4m - 4$$

$$(3) \text{ solved above: } m^3 + m - 2$$

16.2 Computational Tasks

Evaluation

To calculate the probability of a complete configuration, we multiply corresponding conditional probabilities:

$$P(V_1 = x_1, \dots, V_n = x_n) = \prod_{i=1}^n P(V_i = x_i | \mathbf{V}_{\pi(i)} = \mathbf{x}_{\pi(i)})$$

Simulation

For $i = 1, \dots, n$, draw x_j according to $P(V_j = x_j | \mathbf{V}_{\pi(j)} = \mathbf{x}_{\pi(j)})$. Conjoin (x_1, \dots, x_n) to form a complete configuration.

Learning

Learning with predetermined network graph, and from complete observations can be done by direct MLE; i.e., counting.

Inference in Bayesian Networks

One way to handle all inference tasks in Bayesian Networks is to use the brute force method. Using the same summation techniques as in the Joint distribution method, we can do marginalization, conditioning, and completion in Bayesian Networks as well. Although this method does not necessarily use a lot of memory, it is prohibitively expensive in terms of running time. The running time is the same as in the Joint distribution model, i.e., it is exponential. For example, there are $2^5 = 512$ possible configurations in the Burglar-Earthquake example network.

Inference Example using Brute Force

By using the tables we can easily compute the probability of any complete configuration. With appropriate summations and by using the definitions of marginal and conditional probability, we can also solve other inference problems. For example, if we want to calculate $P(B = T | J = T)$, i.e., the probability that a burglar is in the house if John told us over the phone that the alarm is on, we first use the definition of the conditional probability:

$$P(B = T | J = T) = \frac{P(B = T, J = T)}{P(J = T)}$$

The marginal probability $P(B = T, J = T)$ can be calculated using the formula:

$$\begin{aligned} P(B = T, J = T) &= \sum_{E, A, M} P(B = T, E, A, J = T, M) \\ &= \sum_{E, A, M} P(B = T)P(E)P(A|B = T, E)P(J = T|A)P(M|A) \end{aligned}$$

Hence,

$$\begin{aligned} P(B = T, J = T) &= \\ &P(B = T)P(E = T)P(A = T|B = T, E = T)P(J = T|A = T)P(M = T|A = T) \\ &+ P(B = T)P(E = T)P(A = T|B = T, E = T)P(J = T|A = T)P(M = F|A = T) \\ &+ P(B = T)P(E = T)P(A = F|B = T, E = T)P(J = T|A = F)P(M = T|A = F) \\ &+ P(B = T)P(E = T)P(A = F|B = T, E = T)P(J = T|A = F)P(M = F|A = F) \\ &+ P(B = T)P(E = F)P(A = T|B = T, E = F)P(J = T|A = T)P(M = T|A = T) \\ &+ P(B = T)P(E = F)P(A = T|B = T, E = F)P(J = T|A = T)P(M = F|A = T) \\ &+ P(B = T)P(E = F)P(A = F|B = T, E = F)P(J = T|A = F)P(M = T|A = F) \\ &+ P(B = T)P(E = F)P(A = F|B = T, E = F)P(J = T|A = F)P(M = F|A = F) \\ &= 0.001 \cdot 0.002 \cdot 0.95 \cdot 0.9 \cdot 0.7 \\ &+ 0.001 \cdot 0.002 \cdot 0.95 \cdot 0.9 \cdot 0.3 \\ &+ 0.001 \cdot 0.002 \cdot 0.05 \cdot 0.05 \cdot 0.01 \\ &+ 0.001 \cdot 0.002 \cdot 0.05 \cdot 0.05 \cdot 0.99 \\ &+ 0.001 \cdot 0.998 \cdot 0.94 \cdot 0.9 \cdot 0.7 \\ &+ 0.001 \cdot 0.998 \cdot 0.94 \cdot 0.9 \cdot 0.3 \\ &+ 0.001 \cdot 0.998 \cdot 0.06 \cdot 0.05 \cdot 0.01 \\ &+ 0.001 \cdot 0.998 \cdot 0.06 \cdot 0.05 \cdot 0.99 \\ &= 8.49017 \cdot 10^{-4} \end{aligned}$$

To calculate $P(J = T)$, we can represent it as $P(J = T) = P(B = T, J = T) + P(B = F, J = T)$ and first

calculate $P(B = F, J = T)$:

$$\begin{aligned}
 P(B = F, J = T) &= \sum_{E, A, M} P(B = F, E, A, J = T, M) \\
 &= \sum_{E, A, M} P(B = F)P(E)P(A|B = F, E)P(J = T|A)P(M|A) \\
 &= P(B = F)P(E = T)P(A = T|B = F, E = T)P(J = T|A = T)P(M = T|A = T) \\
 &+ P(B = F)P(E = T)P(A = T|B = F, E = T)P(J = T|A = T)P(M = F|A = T) \\
 &+ P(B = F)P(E = T)P(A = F|B = F, E = T)P(J = T|A = F)P(M = T|A = F) \\
 &+ P(B = F)P(E = T)P(A = F|B = F, E = T)P(J = T|A = F)P(M = F|A = F) \\
 &+ P(B = F)P(E = F)P(A = T|B = F, E = F)P(J = T|A = T)P(M = T|A = T) \\
 &+ P(B = F)P(E = F)P(A = T|B = F, E = F)P(J = T|A = T)P(M = F|A = T) \\
 &+ P(B = F)P(E = F)P(A = F|B = F, E = F)P(J = T|A = F)P(M = T|A = F) \\
 &+ P(B = F)P(E = F)P(A = F|B = F, E = F)P(J = T|A = F)P(M = F|A = F) \\
 &= 0.999 \cdot 0.002 \cdot 0.29 \cdot 0.9 \cdot 0.7 \\
 &+ 0.999 \cdot 0.002 \cdot 0.29 \cdot 0.9 \cdot 0.3 \\
 &+ 0.999 \cdot 0.002 \cdot 0.71 \cdot 0.05 \cdot 0.01 \\
 &+ 0.999 \cdot 0.002 \cdot 0.71 \cdot 0.05 \cdot 0.99 \\
 &+ 0.999 \cdot 0.998 \cdot 0.001 \cdot 0.9 \cdot 0.7 \\
 &+ 0.999 \cdot 0.998 \cdot 0.001 \cdot 0.9 \cdot 0.3 \\
 &+ 0.999 \cdot 0.998 \cdot 0.999 \cdot 0.05 \cdot 0.01 \\
 &+ 0.999 \cdot 0.998 \cdot 0.999 \cdot 0.05 \cdot 0.99 \\
 &= 5.12899587 \cdot 10^{-2}
 \end{aligned}$$

Now, we calculate

$$P(J = T) = P(B = T, J = T) + P(B = F, J = T) = 8.49017 \cdot 10^{-4} + 5.12899587 \cdot 10^{-2} = 0.0521389757,$$

and finally

$$P(B = T|J = T) = \frac{P(B = T, J = T)}{P(J = T)} = \frac{8.49017 \cdot 10^{-4}}{0.0521389757} = 0.0162837299467699.$$

Even this small example illustrates inefficiency of this approach.

Complexity of General Inference in Bayesian Networks

In general, we know that for some Bayesian networks the inference must be costly: for example, inference with an arbitrary joint distribution model is expensive, and a joint distribution model can be regarded as a special case of a Bayesian network. The size of the joint distribution model is exponentially large due to the large m^n look up table (m is the size of the domain, and n is the number of random variables).

On the other hand, inference in some very simple Bayesian networks such as the fully independent model and the Naïve Bayes model can be done efficiently. Inference in HMMs can also be done efficiently. An interesting question is: Can we find some other networks where inference is also easy?

We could first notice that the number of parents per node is limited in the special Bayesian networks that we mentioned: in the fully independent model all nodes have 0 parents, in the Naïve Bayes model all nodes have at

most 1 parent, and in the HMM each node has at most 1 parent as well. So, the first idea is to restrict the maximal number of parents in each node.

If the maximal number of parents is 1, then inference can be done efficiently, and we will see this soon. If the maximal number of parents is 2, inference can be hard. The representation size (the table size) is still polynomial, but the inference task can be NP-hard. We can prove that inference in Bayesian Networks with at most 2 parents is NP-hard by the NP-reduction from the Circuit satisfiability problem, or, equivalently, to satisfiability of an arbitrary boolean function. A sketch of this proof is as follows: If we have an arbitrary circuit satisfiability problem it can be represented as a directional circuit with some 0/1 inputs, boolean gates (AND, OR, and NOT), and an output. The satisfiability problem is the question of whether it is possible to set input nodes to the values 0 or 1 in such a way that we obtain 1 at the output node. We can turn this circuit into a Bayesian Network but attaching a node to each input and assigning probability 0.5 to 0 or 1 values. The probability tables are simple 0-1 tables corresponding to associated boolean operations. Satisfiability of the circuit is then equivalent to verifying that the marginal probability of the output node is greater than zero.

We will show that a Bayesian network in which each node has at most one parent node allows for efficient inference. Actually, more generally, if the Bayesian network graph has a shape of the tree; i.e., its undirected version of the graph is a tree, then we know how to do inference efficiently.

16.3 Sum-Product Algorithms for Bayesian Networks

Slide notes:

Sum-Product Algorithms for Bayesian Networks

- Basic idea: optimizing sum-product calculation using graph structure
Described in “Factor graphs and the Sum-Product Algorithm” by Kschishang, Frey, and Loeliger in 2000
- Algorithm overview:
 1. Construction of a factor graph
 2. Message-passing algorithms
- Construction of the factor graph
- Principles of message passing

The algorithms for inference (marginalization, conditioning, and completion) in tree Bayesian networks can be formulated conveniently as message passing algorithms, also known as product-sum algorithms. The algorithm was re-invented in different disguises in Hidden Markov Models, Kalman filters, Error-correcting codes, Bayesian Networks, and similar.²

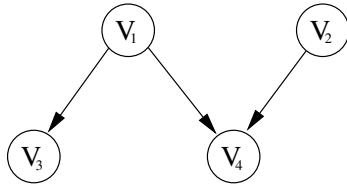
Before presenting the message passing algorithms, let us first see on an example how we can do marginalization efficiently on a Bayesian network, using the BN assumption and the sum separation rule:

Example

Note: This example is not covered in class.

This example illustrates the underlying mathematical computation that leads to the message passing algorithms that we will cover. Let us assume that we have the following Bayesian Network:

²The paper introducing the message passing algorithms is “Factor Graphs and the Sum-Product Algorithm” by Frank R. Kschishang, Brendan J. Frey, and Hans-Andrea Loeliger in 2000.



and that we want to compute $P(V_4 = x_4)$:

$$\begin{aligned}
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} P(V_1 = x_1, V_2 = x_2, V_3 = x_3, V_4 = x_4) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} P(V_1 = x_1) P(V_2 = x_2) P(V_3 = x_3 | V_1 = x_1) \\
 &\quad P(V_4 = x_4 | V_1 = x_1, V_2 = x_2)
 \end{aligned}$$

This computation has $m^3 - 1$ additions and $3m^3$ multiplications, which is very expensive (as expensive as the joint distribution model).

Let us think of conditional probabilities as functions:

$$\begin{aligned}
 f_1(x_1) &= P(V_1 = x_1) \\
 f_2(x_2) &= P(V_2 = x_2) \\
 f_3(x_1, x_3) &= P(V_3 = x_3 | V_1 = x_1) \\
 f_4(x_1, x_2, x_4) &= P(V_4 = x_4 | V_1 = x_1, V_2 = x_2)
 \end{aligned}$$

and let us try to calculate the value of the above expression with less additions and multiplications. We use the fact that summations of factors which depend on variables x_2 and x_3 can be separated.

$$\begin{aligned}
 P(V_4 = x_4) &= \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} f_1(x_1) f_2(x_2) f_3(x_1, x_3) f_4(x_1, x_2, x_4) \\
 &= \sum_{x_1} f_1(x_1) \sum_{x_3} f_3(x_1, x_3) \sum_{x_2} f_2(x_2) f_4(x_1, x_2, x_4) \\
 &= \sum_{x_1} f_1(x_1) \left(\sum_{x_2} f_2(x_2) f_4(x_1, x_2, x_4) \right) \sum_{x_3} f_3(x_1, x_3)
 \end{aligned}$$

We can calculate number of operations: First, for each value of x_1 we make m multiplications in sum over x_2 , $m - 1$ additions in sum over x_2 , $m - 1$ additions in sum over x_3 , and two more multiplications of these two sums and f_1 , which adds up to $m + 2$ multiplications and $2m - 2$ additions. This is repeated m times, which makes $m^2 + 2m$ multiplications and $2m^2 - 2m$ additions. After that, we make additional $m - 1$ additions, which makes $m^2 + 2m$ multiplications, and $2m^2 - m - 1$ additions.

This is significantly better than $3m^3$ multiplications and $m^3 - 1$ additions.

Breaking up a brute-force formula as in this example into more efficient computations gets more complicated and harder to manage in general Bayesian Networks. It can be done relatively efficiently in tree-structured Bayesian Networks; i.e., networks that do not have cycles when they graphs are treated as undirected graphs.

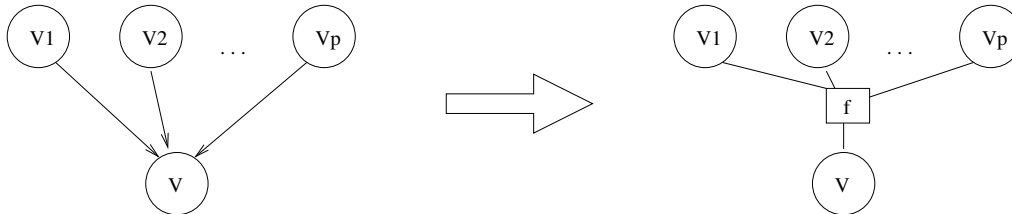
Message Passing Algorithms Framework

Applying the message passing algorithms to Bayesian Network inference consists of two steps:

1. translation of a BN into a factor graph, and
2. message exchange between nodes in the factor graph.

Factor graph

The first step in Bayesian network inference is to convert the Bayesian network graph into a **factor graph**. A factor graph is an undirected graph in which each variable is represented with one “round” node and one “square” node representing the corresponding conditional probability table (CPT). The square node is connected with all variables appearing in CPT. So, if V is a variable in a Bayesian network, the following transformation is applied.

**Example:**

Using the Burglar-Earthquake example:

