

CSCI 4152/6509 — Natural Language Processing

26-Oct-2009

Lecture 18: Bayesian Belief Networks

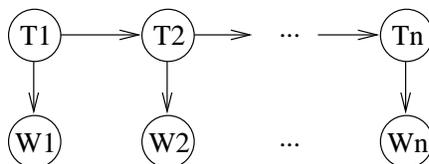
Room: FASS 2176
Time: 11:35 – 12:25

Previous Lecture

- Smoothing:
 - Add-one (Laplace) smoothing,
 - Bell-Witten smoothing;
- Hidden Markov Model,
 - graphical representations,
 - assumption,
 - POS example

11.1 HMM POS Example

When using HMMs for POS tagging, we assume that the hidden internal states of the HMM correspond to correct POS tags of words, while the words correspond to generated observed variables. According to this, a sentence of n words would be associated with the following HMM graph:



The variables W_1, \dots, W_n are assigned to words in the sentence, while variables T_1, \dots, T_n are assigned POS tags.

All tables of conditional probabilities for $P(T_2|T_1)$, $P(T_3|T_2)$, \dots , $P(T_n|T_{n-1})$, as well as tables $P(W_1|T_1)$, $P(W_2|T_2)$, \dots , $P(W_n|T_n)$, are equal.

Having this in mind, suppose that we are given the following training data:

```
swat V flies N like P ants N
time N flies V like P an D arrow N
```

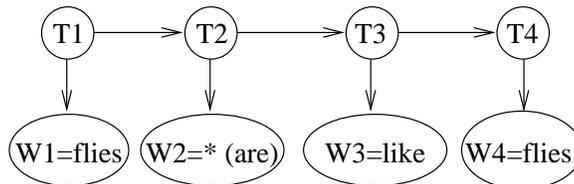
To accommodate for unseen words, we can assign a special symbol $*$ to unknown words, and assume that it occurred 0.5 “times” with each tag.

The following probability tables are generated using this smoothing technique:

T_1	$P(T_1)$	T_{i-1}	T_i	$P(T_i T_{i-1})$	and	T_i	W_i	$P(W_i T_i)$
N	0.5	D	N	1		D	an	$2/3 \approx 0.666666667$
V	0.5	N	P	0.5		D	*	$1/3 \approx 0.333333333$
		N	V	0.5		N	ants	$2/9 \approx 0.222222222$
		P	D	0.5		N	arrow	$2/9 \approx 0.222222222$
		P	N	0.5		N	flies	$2/9 \approx 0.222222222$
		V	N	0.5		N	time	$2/9 \approx 0.222222222$
		V	P	0.5		N	*	$1/9 \approx 0.111111111$
						P	like	0.8
						P	*	0.2
						V	flies	0.4
						V	swat	0.4
						V	*	0.2

The above tables may not seem to be complete in the sense that some combinations of variable assignments are not included. For example, $P(T_i = V|T_{i-1} = D)$ is not included. For the sake of saving space, we have not included the combinations having probability 0, such as this one.

Let us use the Hidden Markov Model to POS tag the sentence “flies are like flies.”



The problem of POS tagging in this case is the problem of finding the most probable values of the variables T_i given the values of variables W_i , i.e., it is the completion problem of finding

$$\begin{aligned}
 \arg \max_T P(T|W = \text{sentence}) &= \arg \max_T \frac{P(T, W = \text{sentence})}{P(W = \text{sentence})} = \arg \max_T P(T, W = \text{sentence}) \\
 &= \arg \max_T P(T_1) \cdot P(W_1 = \text{flies}|T_1) \cdot P(T_2|T_1) \cdot P(W_2 = *|T_2) \\
 &\quad \cdot P(T_3|T_2) \cdot P(W_3 = \text{like}|T_3) \cdot P(T_4|T_3) \cdot P(W_4 = \text{flies}|T_4)
 \end{aligned}$$

where T and W denote arrays of variables T_i and W_i . One way to find the values of T that maximize the given probability is to test all variations of their values. This number is exponential in general, and in this case it is $4^4 = 256$. A much more efficient solution can be obtained by applying a dynamic programming approach, known as the Viterbi algorithm. The idea of the Viterbi algorithm is to incrementally calculate maximal values of the following parts of the above product:

$$P(T_1) \cdot P(W_1 = \text{flies}|T_1)$$

for all possible values of T_1 , then

$$P(T_1) \cdot P(W_1 = \text{flies}|T_1) \cdot P(T_2|T_1) \cdot P(W_2 = *|T_2)$$

for all possible values of T_2 and so on. The computation can be summarized in the following table:

	$T_1 (W_1 = \text{flies})$	$T_2 (W_2 = *)$	$T_3 (W_3 = \text{like})$	$T_4 (W_4 = \text{flies})$
	$P(T_1)P(W_1 T_1)$	$p \cdot P(T_2 T_1)P(W_2 T_2)$	$p \cdot P(T_2 T_1)P(W_2 T_2)$	$p \cdot P(T_2 T_1)P(W_2 T_2)$
D	$0 \times 0 = 0$	DD: $0 \times 0 \times \frac{1}{3} = 0$ ND: $\frac{1}{9} \times 0 \times \frac{1}{3} = 0$ PD: 0 VD: 0 max: 0	DD: $0 \times 0 \times 0 = 0$ ND: $\frac{1}{90} \times 0 \times 0 = 0$ PD: $\frac{1}{50} \times \frac{1}{2} \times 0 = 0$ VD: $\frac{1}{90} \times 0 \times 0 = 0$ max: 0	DD: $0 \times 0 \times 0 = 0$ ND: $0 \times 0 \times 0 = 0$ PD: $\frac{1}{225} \times 0.5 \times 0 = 0$ VD: $0 \times 0 \times 0 = 0$ max: 0
N	$0.5 \times \frac{2}{9} = \frac{1}{9}$	DN: $0 \times 1 \dots = 0$ NN: $\frac{1}{9} \times 0 \dots = 0$ PN: $0 \times \dots = 0$ VN: $0.2 \times 0.5 \times \frac{1}{9} = \frac{1}{90}$ max: $\frac{1}{90}$	DN: $0 \times 1 \times 0 = 0$ NN: $\frac{1}{90} \times 0 \dots = 0$ PN: $\frac{1}{50} \times 0.5 \times 0 = 0$ VN: $\frac{1}{90} \times 0.5 \times 0 = 0$ max: 0	DN: $0 \times 1 \times \frac{2}{9} = 0$ NN: $0 \times 0 \times \frac{2}{9} = 0$ PN: $\frac{1}{225} \times 0.5 \times \frac{2}{9} = \frac{1}{2025}$ VN: $0 \times 0.5 \times \frac{2}{9} = 0$ max: $\frac{1}{2025}$
P	$0 \times 0 = 0$	DP: $0 \times \dots = 0$ NP: $\frac{1}{9} \times 0.5 \times 0.2 = \frac{1}{90}$ PP: $0 \times \dots = 0$ VP: $0.2 \times 0.5 \times 0.2 = \frac{1}{50}$ max: $\frac{1}{50}$	DP: $0 \times 0 \times 0.8 = 0$ NP: $\frac{1}{90} \times 0.5 \times 0.8 = \frac{1}{225}$ PP: $\frac{1}{50} \times 0 \times 0.8 = 0$ VP: $\frac{1}{90} \times 0.5 \times 0.8 = \frac{1}{225}$ max: $\frac{1}{225}$	DP: $0 \times 0 \times 0 = 0$ NP: $0 \times 0.5 \times 0 = 0$ PP: $\frac{1}{225} \times 0 \times 0 = 0$ VP: $0 \times 0.5 \times 0 = 0$ max: 0
V	$0.5 \times 0.4 = 0.2$	DV: $0 \times \dots = 0$ NV: $\frac{1}{9} \times 0.5 \times 0.2 = \frac{1}{90}$ PV: $0 \times \dots = 0$ VV: $0.2 \times 0 \dots = 0$ max: $\frac{1}{90}$	DV: $0 \times 0 \times 0 = 0$ NV: $\frac{1}{90} \times 0.5 \times 0 = 0$ PV: $\frac{1}{50} \times 0 \times 0 = 0$ VV: $\frac{1}{90} \times 0 \times 0 = 0$ max: 0	DV: $0 \times 0 \times 0.4 = 0$ NV: $0 \times 0.5 \times 0.4 = 0$ PV: $\frac{1}{225} \times 0 \times 0.4 = 0$ VV: $0 \times 0 \times 0.4 = 0$ max: 0

The table is filled column by column. We can see now that the largest value that the expression

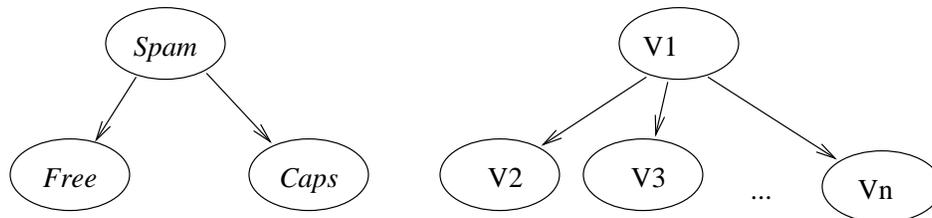
$$P(T_1) \cdot P(W_1 = \text{flies}|T_1) \cdot P(T_2|T_1) \cdot P(W_2 = *|T_2) \cdot P(T_3|T_2) \cdot P(W_3 = \text{like}|T_3) \cdot P(T_4|T_3) \cdot P(W_4 = \text{flies}|T_4)$$

can obtain is $\frac{1}{2025}$, and is achieved with $T_4 = N$. If we work backwards through the table, we can obtain the optimal values for previous variables as well: $T_3 = P$, $T_2 = V$, and $T_1 = N$. We can also choose $T_2 = N$, but in this case we have $T_1 = V$.

12 Bayesian networks

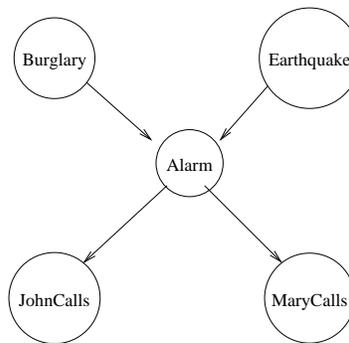
Bayesian Networks (also known as **belief networks**, **Bayesian belief networks**, or **decision networks**) provide a way to create structured probabilistic models. They make a balance between too strong independence assumptions, as the ones found in Naïve Bayes model or Hidden Markov Model, and the number of parameters of the joint distribution model.

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



Examples:

Let us 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	$P(B)$	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
				F	T	T	0.94	F	T	0.05	F	T	0.01
				F	T	F	0.06	F	F	0.95	F	F	0.99
				F	F	T	0.29						
				F	F	F	0.71						
							0.001						
							0.999						

– Inference by 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.

After this informal description of Bayesian Network, here is a more formal definition:

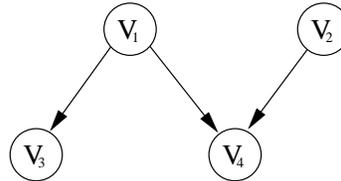
Definition 12.1 (Bayesian Network) A Bayesian Network is defined by a directed acyclic graph (DAG) and a collection of conditional probability tables, where nodes in the graph represent random variables and directed edges in the graph represent conditional independence assumptions. The edges are interpreted in the following way: If V_j ($1 \leq j \leq n$) is a random variable, and $\mathbf{V}_{\pi(j)}$ are parent variables of V_j , i.e., all source nodes for edges whose destination node is V_j , then the probability of V_j given variables $\mathbf{V}_{\pi(j)}$ is independent of any other variable; i.e.,

$$P(V_j = x_j | \text{all other variables}) = P(V_j = x_j | \mathbf{V}_{\pi(j)} = \mathbf{x}_{\pi(j)})$$

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 $P(V_i | \mathbf{V}_{\pi(i)})$. 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

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

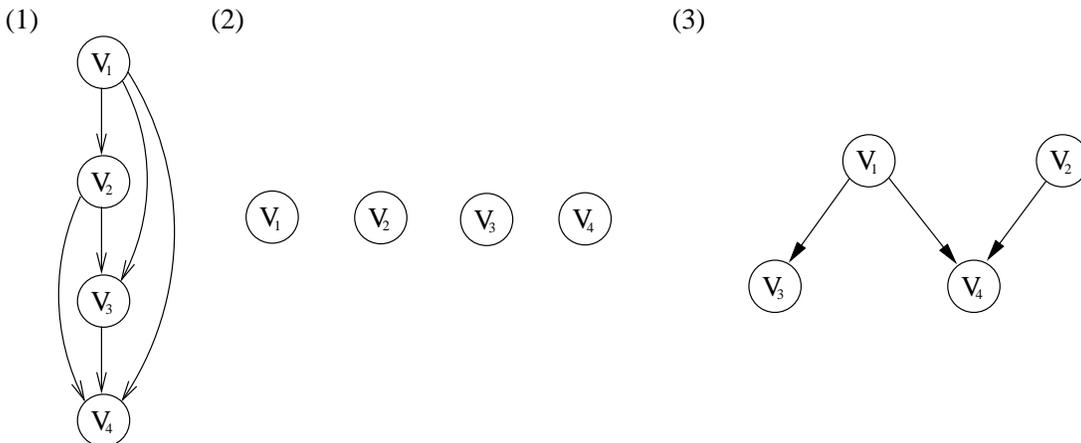
For the above network, the number of free parameters is:

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

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) $(m - 1) + (m^2 - m) + (m^3 - m^2) + (m^4 - m^3) = m^4 - 1$
- (2) $(m - 1) + (m - 1) + (m - 1) + (m - 1) = 4m - 4$
- (3) solved above: $m^3 + m - 2$

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
				F	T	T	0.94	F	T	0.05	F	T	0.01
				F	T	F	0.06	F	F	0.95	F	F	0.99
				F	F	T	0.29	F					
				F	F	F	0.71	F					
				F	F	F	0.001						
				F	F	F	0.999						

12.1 Computational Tasks**Evaluation**

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

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

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.