

Artificial Intelligence

Bayesian networks
Chapter 14, AIMA

Inference

- Inference in the statistical setting means computing probabilities for different outcomes to be true given the information

$$\mathbf{P}(\textit{Outcome} \mid \textit{Information})$$

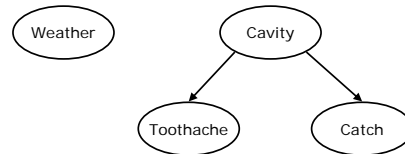
- We need an efficient method for doing this, which is more powerful than the naïve Bayes model.

Bayesian networks

A **Bayesian network** is a directed graph in which each node is annotated with quantitative probability information:

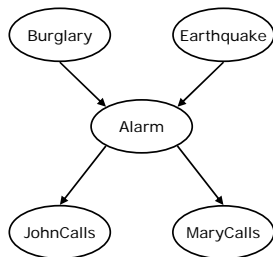
1. A set of random variables, $\{X_1, X_2, X_3, \dots\}$, makes up the nodes of the network.
2. A set of directed links connect pairs of nodes, parent \rightarrow child
3. Each node X_i has a conditional probability distribution $\mathbf{P}(X_i \mid \textit{Parents}(X_i))$.
4. The graph is a directed acyclic graph (DAG).

The dentist network



	Toothache		¬Toothache	
Cavity	catch	¬catch	catch	¬catch
catch	0.1000	0.0200	0.0700	0.0000
¬catch	0.0300	0.2000	0.1500	0.2500

The alarm network



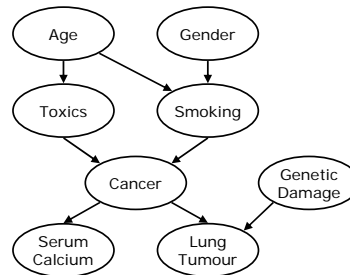
Burglar alarm responds to both earthquakes and burglars.

Two neighbors: John and Mary, who have promised to call you when the alarm goes off.

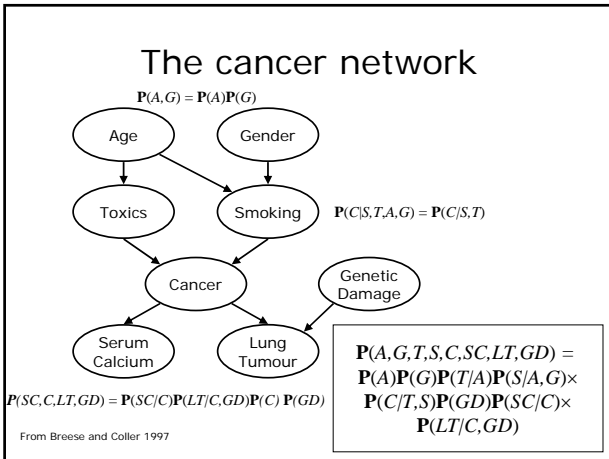
John always calls when there's an alarm, and sometimes when there's not an alarm.

Mary sometimes misses the alarms (she likes loud music).

The cancer network



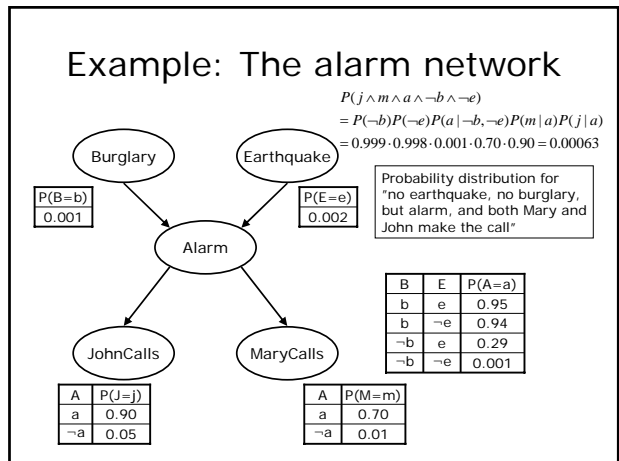
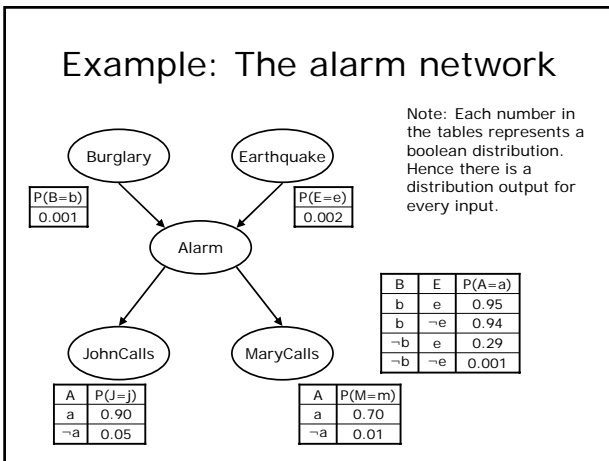
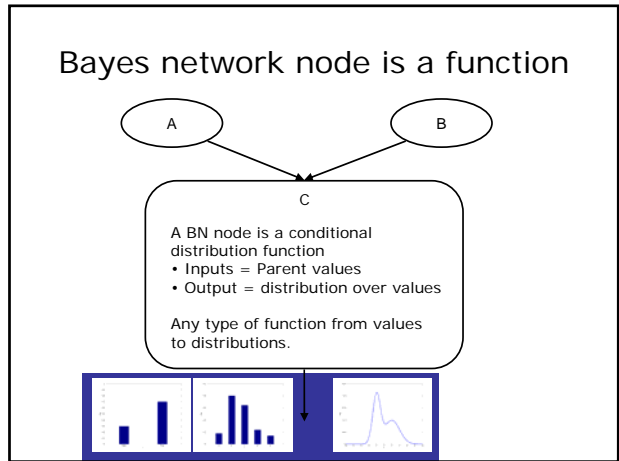
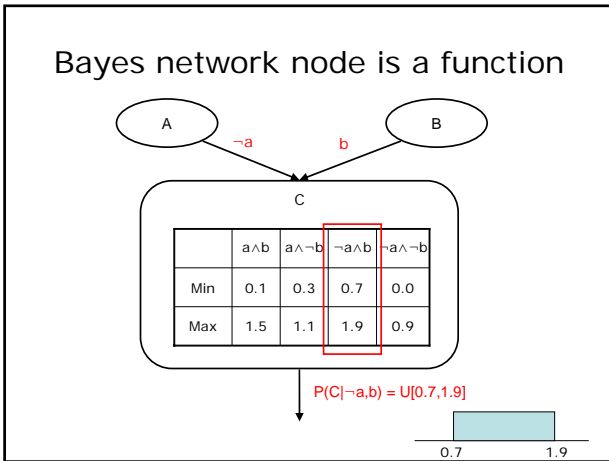
From Breese and Coller 1997



The product (chain) rule

$$P(X_1 = x_1 \wedge X_2 = x_2 \wedge \dots \wedge X_n = x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

(This is for Bayesian networks, the general case comes later in this lecture)



Meaning of Bayesian network

The general chain rule (always true):

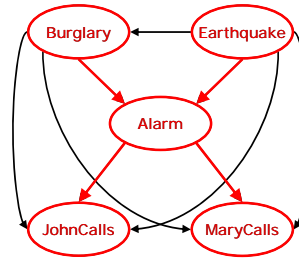
$$P(x_1, x_2, \dots, x_n) = P(x_1 | x_2, x_3, \dots, x_n) P(x_2, x_3, \dots, x_n) = \\ P(x_1 | x_2, x_3, \dots, x_n) P(x_2 | x_3, x_4, \dots, x_n) P(x_3, x_4, \dots, x_n) = \dots \\ = \prod_{i=1}^n P(x_i | x_{i+1}, \dots, x_n)$$

The Bayesian network chain rule:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

The BN is a correct representation of the domain iff each node is conditionally independent of its predecessors, given its parents.

The alarm network



The fully correct alarm network might look something like the figure.

The Bayesian network (red) assumes that some of the variables are independent (or that the dependencies can be neglected since they are very weak).

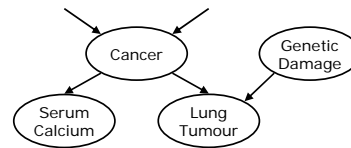
The correctness of the Bayesian network of course depends on the validity of these assumptions.

It is this sparse connection structure that makes the BN approach feasible (~linear growth in complexity rather than exponential)

How construct a BN?

- Add nodes in causal order ("causal" determined from expertise).
- Determine conditional independence using either (or all) of the following semantics:
 - Blocking/d-separation rule
 - Non-descendant rule
 - Markov blanket rule
 - Experience/your beliefs

Path blocking & d-separation



Intuitively, knowledge about Serum Calcium influences our belief about Cancer, if we don't know the value of Cancer, which in turn influences our belief about Lung Tumour, etc.

However, if we are given the value of Cancer (i.e. $C = \text{true}$ or false), then knowledge of Serum Calcium will not tell us anything about Lung Tumour that we don't already know.

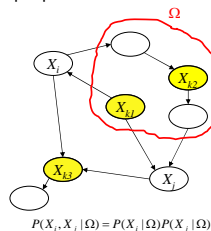
We say that Cancer **d-separates** (direction-dependent separation) Serum Calcium and Lung Tumour.

Path blocking & d-separation

X_i and X_j are d-separated if all paths between them are blocked

Two nodes X_i and X_j are conditionally independent given a set $\Omega = \{X_1, X_2, X_3, \dots\}$ of nodes if for every undirected path in the BN between X_i and X_j there is some node X_k on the path having one of the following three properties:

1. $X_k \in \Omega$, and both arcs on the path lead out of X_k .
2. $X_k \in \Omega$, and one arc on the path leads into X_k and one arc leads out.
3. Neither X_k nor any descendant of X_k is in Ω , and both arcs on the path lead into X_k .

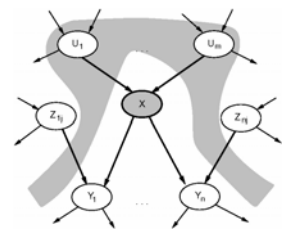


X_k **blocks** the path between X_i and X_j

$$P(X_i, X_j | \Omega) = P(X_i | \Omega) P(X_j | \Omega)$$

Non-descendants

A node is conditionally independent of its non-descendants (Z_{ij}), given its parents.

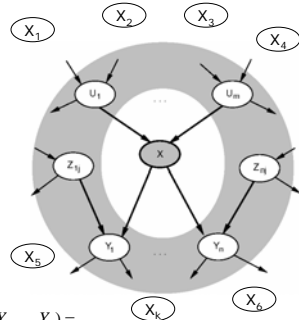


$$P(X, Z_{1j}, \dots, Z_{nj} | U_1, \dots, U_m) =$$

$$P(X | U_1, \dots, U_m) P(Z_{1j}, \dots, Z_{nj} | U_1, \dots, U_m)$$

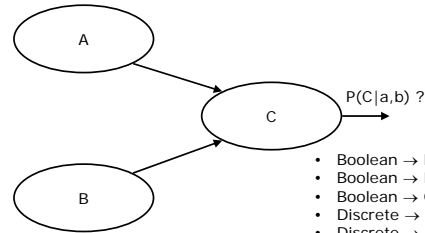
Markov blanket

A node is conditionally independent of all other nodes in the network, given its parents, children, and children's parents. These constitute the nodes *Markov blanket*.



$$P(X, X_1, \dots, X_k | U_1, \dots, U_m, Z_{1j}, \dots, Z_{nj}, Y_1, \dots, Y_n) = P(X | U_1, \dots, U_m, Z_{1j}, \dots, Z_{nj}, Y_1, \dots, Y_n) P(X_1, \dots, X_k | U_1, \dots, U_m, Z_{1j}, \dots, Z_{nj}, Y_1, \dots, Y_n)$$

Efficient representation of PDs



- Boolean → Boolean
- Boolean → Discrete
- Boolean → Continuous
- Discrete → Boolean
- Discrete → Discrete
- Discrete → Continuous
- Continuous → Boolean
- Continuous → Discrete
- Continuous → Continuous

Efficient representation of PDs

Boolean → Boolean:

Noisy-OR, Noisy-AND

Boolean/Discrete → Discrete:

Noisy-MAX

Bool./Discr./Cont. → Continuous:

Parametric distribution (e.g. Gaussian)

Continuous → Boolean:

Logit/Probit

Noisy-OR example

Boolean → Boolean

		P(E C ₁ ,C ₂ ,C ₃)							
C ₁		0	1	0	0	1	1	0	1
C ₂		0	0	1	0	1	0	1	1
C ₃		0	0	0	1	0	1	1	1
P(E=0)		1	0.1	0.1	0.1	0.01	0.01	0.01	0.001
P(E=1)		0	0.9	0.9	0.9	0.99	0.99	0.99	0.999

The effect (E) is off (false) when none of the causes are true. The probability for the effect increases with the number of true causes.

$$P(E=0) = 10^{-(\#True)} \quad (\text{for this example})$$

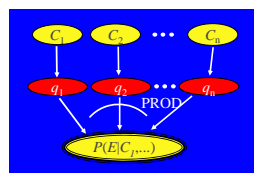
Example from L.E. Sucar

Noisy-OR general case

Boolean → Boolean

$$P(E=0 | C_1, C_2, \dots, C_n) = \prod_{i=1}^n q_i^{C_i}$$

$$C_i = \begin{cases} 1 & \text{if true} \\ 0 & \text{if false} \end{cases}$$



Example on previous slide used
 $q_i = 0.1$ for all i .

Image adapted from Laskey & Mahoney 1999

Noisy-MAX

Boolean → Discrete

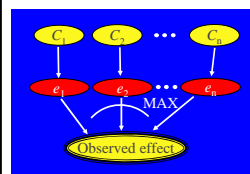


Image adapted from Laskey & Mahoney 1999

Effect takes on the max value from different causes

Restrictions:

- Each cause must have an off state, which does not contribute to effect
- Effect is off when all causes are off
- Effect must have consecutive escalating values: e.g., absent, mild, moderate, severe.

$$P(E = e_k | C_1, C_2, \dots, C_n) = \prod_{i=1}^n q_{i,k}^{C_i}$$

Parametric probability densities

Boolean/Discr./Continuous → Continuous

Use parametric probability densities, e.g., the normal distribution

$$P(X) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = N(\mu, \sigma)$$

Gaussian networks (a = input to the node)

$$P(X) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-\alpha-\beta a)^2}{2\sigma^2}\right]$$

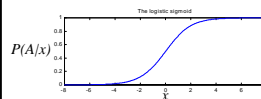
Probit & Logit

Discrete → Boolean

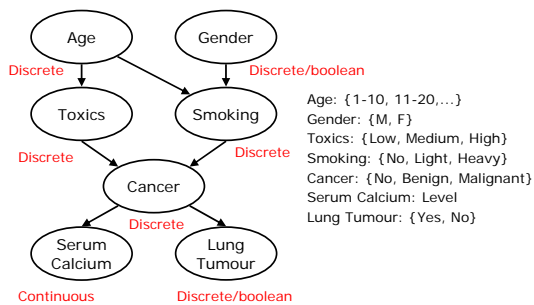
If the input is continuous but output is boolean, use probit or logit

$$\text{Logit: } P(A = a | x) = \frac{1}{1 + \exp[-2(\mu - x) / \sigma]}$$

$$\text{Probit: } P(A = a | x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp(-(x-\mu)^2 / \sigma^2) dx$$



The cancer network



Inference in BN

Inference means computing $\mathbf{P}(X|\mathbf{e})$, where X is a query (variable) and \mathbf{e} is a set of evidence variables (for which we know the values).

Examples:

$\mathbf{P}(\text{Burglary} | \text{john_calls}, \text{mary_calls})$

$\mathbf{P}(\text{Cancer} | \text{age}, \text{gender}, \text{smoking}, \text{serum_calcium})$

$\mathbf{P}(\text{Cavity} | \text{toothache}, \text{catch})$

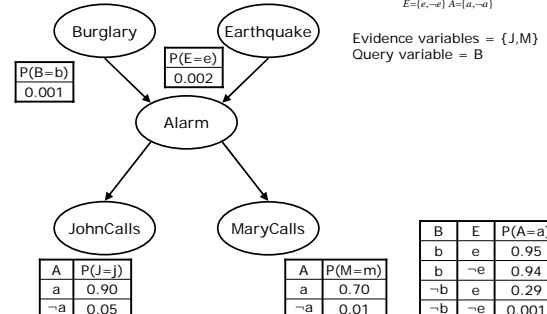
Exact inference in BN

$$\mathbf{P}(X | \mathbf{e}) = \frac{\mathbf{P}(X, \mathbf{e})}{\mathbf{P}(\mathbf{e})} = \alpha \mathbf{P}(X, \mathbf{e}) = \alpha \sum_{\mathbf{y}} \mathbf{P}(X, \mathbf{e}, \mathbf{y})$$

"Doable" for boolean variables: Look up entries in conditional probability tables (CPTs).

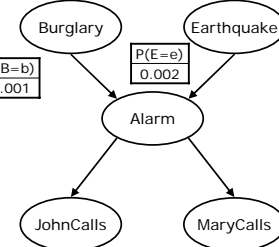
Example: The alarm network

What is the probability for a burglary if both John and Mary call? $\mathbf{P}(B | j, m) = \alpha \sum_{E=\{e, \neg e\}} \sum_{A=\{a, \neg a\}} \mathbf{P}(B, E, A, j, m)$



Example: The alarm network

What is the probability for a burglary if both John and Mary call? $P(B|j,m) = \alpha \sum_{E=\{e,-e\}} \sum_{A=\{a,-a\}} P(B,E,A,j,m)$



$$P(b,j,m) = P(j|b,E,A)P(b,E,A) = P(j|A)P(m|A)P(a|b,E)P(b,E) = P(j|A)P(m|A)P(a|b,E)P(b)P(E) = 10^{-3} \times P(j|A)P(m|A)P(A|b,E)P(E)$$

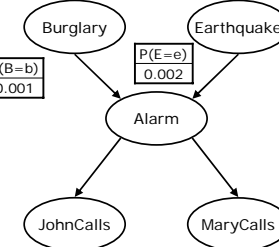
A	P(J=j)
a	0.90
-a	0.05

A	P(M=m)
a	0.70
-a	0.01

B	E	P(A=a)
b	e	0.95
b	-e	0.94
-b	e	0.29
-b	-e	0.001

Example: The alarm network

What is the probability for a burglary if both John and Mary call? $P(B|j,m) = \alpha \sum_{E=\{e,-e\}} \sum_{A=\{a,-a\}} P(B,E,A,j,m)$



$$P(b,j,m) = 10^{-3} \sum_{A=\{a,-a\}} P(j|A)P(m|A)P(A|b,E)P(E) = 10^{-3} [P(j|a)P(m|a)P(a|b,e)P(e) + P(j|a)P(m|a)P(a|b,-e)P(-e) + P(j|-a)P(m|-a)P(-a|b,e)P(e) + P(j|-a)P(m|-a)P(-a|b,-e)P(-e)] = 0.5923 \times 10^{-3}$$

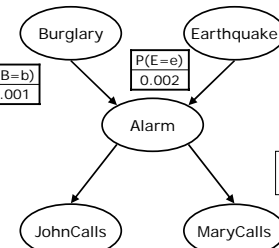
A	P(J=j)
a	0.90
-a	0.05

A	P(M=m)
a	0.70
-a	0.01

B	E	P(A=a)
b	e	0.95
b	-e	0.94
-b	e	0.29
-b	-e	0.001

Example: The alarm network

What is the probability for a burglary if both John and Mary call? $P(B|j,m) = \alpha \sum_{E=\{e,-e\}} \sum_{A=\{a,-a\}} P(B,E,A,j,m)$



$$P(b,j,m) = 0.5923 \times 10^{-3}$$

$$P(-b,j,m) = 1.491 \times 10^{-3}$$

$$\alpha = P(j,m)^{-1} = [P(b,j,m) + P(-b,j,m)]^{-1} = [2.083 \times 10^{-3}]^{-1}$$

$$P(b|j,m) = \alpha P(b,j,m) = 0.284$$

$$P(-b|j,m) = \alpha P(-b,j,m) = 0.716$$

A	P(J=j)
a	0.90
-a	0.05

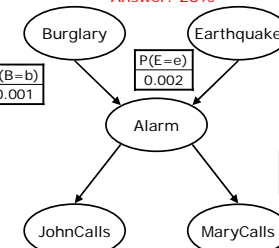
A	P(M=m)
a	0.70
-a	0.01

B	E	P(A=a)
b	e	0.95
b	-e	0.94
-b	e	0.29
-b	-e	0.001

Example: The alarm network

What is the probability for a burglary if both John and Mary call? $P(B|j,m) = \alpha \sum_{E=\{e,-e\}} \sum_{A=\{a,-a\}} P(B,E,A,j,m)$

Answer: 28%



$$P(b,j,m) = 0.5923 \times 10^{-3}$$

$$P(-b,j,m) = 1.491 \times 10^{-3}$$

$$\alpha = P(j,m)^{-1} = [P(b,j,m) + P(-b,j,m)]^{-1} = [2.083 \times 10^{-3}]^{-1}$$

$$P(b|j,m) = \alpha P(b,j,m) = 0.284$$

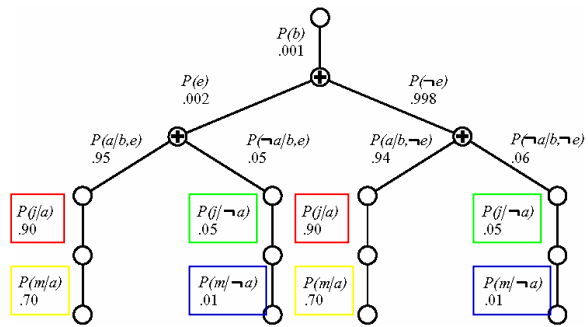
$$P(-b|j,m) = \alpha P(-b,j,m) = 0.716$$

A	P(J=j)
a	0.90
-a	0.05

A	P(M=m)
a	0.70
-a	0.01

B	E	P(A=a)
b	e	0.95
b	-e	0.94
-b	e	0.29
-b	-e	0.001

Use depth-first search



A lot of unnecessary repeated computation...

Complexity of exact inference

- By eliminating repeated calculation & uninteresting paths we can speed up the inference a lot.
- Linear time complexity for singly connected networks (polytrees).
- Exponential for multiply connected networks.
 - Clustering can improve this

Approximate inference in BN

- Exact inference is intractable in large multiply connected BNs \Rightarrow use approximate inference: Monte Carlo methods (random sampling).
 - Direct sampling
 - Rejection sampling
 - Likelihood weighting
 - **Markov chain Monte Carlo**

Markov chain Monte Carlo

1. Fix the evidence variables (E_1, E_2, \dots) at their given values.
2. Initialize the network with values for all other variables, including the query variable.
3. Repeat the following many, many, many times:
 - a. Pick a non-evidence variable at random (query X_i or hidden Y_j)
 - b. Select a new value for this variable, conditioned on the current values in the variable's Markov blanket.

Monitor the values of the query variables.