

Tutorial on Bayesian Networks

Jack Breese

Microsoft Research

breese@microsoft.com

Daphne Koller

Stanford University

koller@cs.stanford.edu

First given as a AAAI'97 tutorial.

Probabilities

- Probability distribution $P(X|\xi)$
 - ◆ X is a random variable
 - Discrete
 - Continuous
 - ◆ ξ is background state of information

Discrete Random Variables

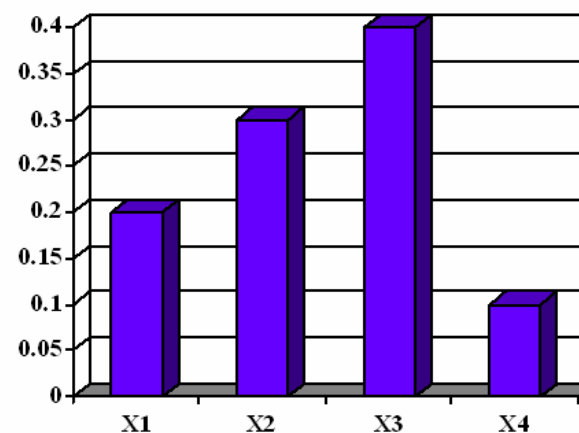
- Finite set of possible outcomes

$$X \in \{x_1, x_2, x_3, \dots, x_n\}$$

$$P(x_i) \geq 0$$

$$\sum_{i=1}^n P(x_i) = 1$$

$$X \text{ binary: } P(x) + P(\bar{x}) = 1$$



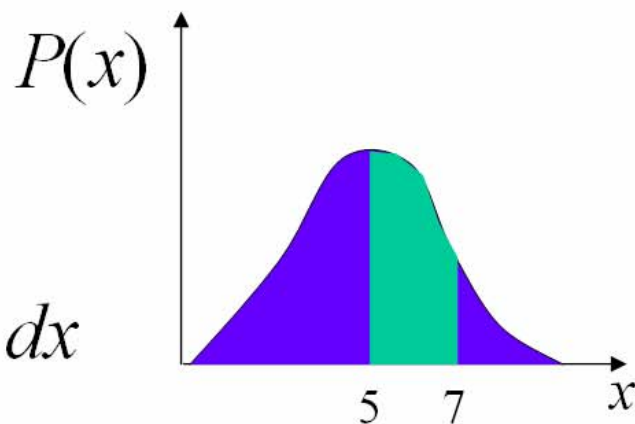
Continuous Random Variable

- Probability distribution (density function) over continuous values

$$X \in [0,10] \quad P(x) \geq 0$$

$$\int_0^{10} P(x) dx = 1$$

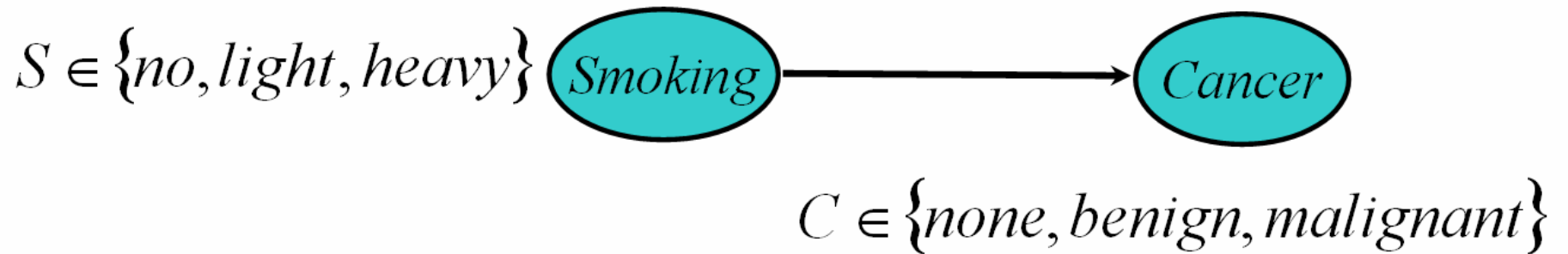
$$P(5 \leq x \leq 7) = \int_5^7 P(x) dx$$



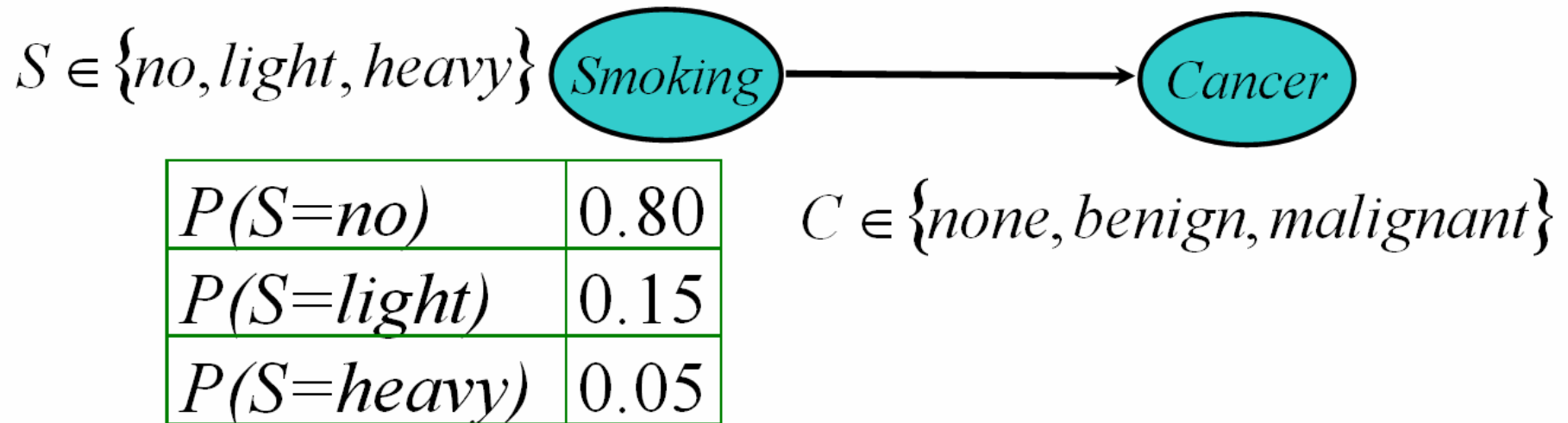
Bayesian networks

- Basics
 - ◆ Structured representation
 - ◆ Conditional independence
 - ◆ Naïve Bayes model
 - ◆ Independence facts

Bayesian Networks



Bayesian Networks



Bayesian Networks



$P(S=no)$	0.80
$P(S=light)$	0.15
$P(S=heavy)$	0.05

$C \in \{none, benign, malignant\}$

$Smoking=$	no	$light$	$heavy$
$P(C=none)$	0.96	0.88	0.60
$P(C=benign)$	0.03	0.08	0.25
$P(C=malig)$	0.01	0.04	0.15

Product Rule

■ $P(C,S) = P(C|S) P(S)$

$S \Downarrow$ $C \Rightarrow$	<i>none</i>	<i>benign</i>	<i>malignant</i>
<i>no</i>	0.768	0.024	0.008
<i>light</i>	0.132	0.012	0.006
<i>heavy</i>	0.035	0.010	0.005

Marginalization

$S \downarrow$ $C \Rightarrow$	<i>none</i>	<i>benign</i>	<i>malig</i>	total
<i>no</i>	0.768	0.024	0.008	.80
<i>light</i>	0.132	0.012	0.006	.15
<i>heavy</i>	0.035	0.010	0.005	.05
total	0.935	0.046	0.019	

$P(\text{Smoke})$

$P(\text{Cancer})$

Bayes Rule Revisited

$$P(S|C) = \frac{P(C|S)P(S)}{P(C)} = \frac{P(C,S)}{P(C)}$$

Bayes Rule Revisited

$$P(S|C) = \frac{P(C|S)P(S)}{P(C)} = \frac{P(C,S)}{P(C)}$$

$S \Downarrow$ $C \Rightarrow$	<i>none</i>	<i>benign</i>	<i>malig</i>
<i>no</i>	0.768/.935	0.024/.046	0.008/.019
<i>light</i>	0.132/.935	0.012/.046	0.006/.019
<i>heavy</i>	0.030/.935	0.015/.046	0.005/.019

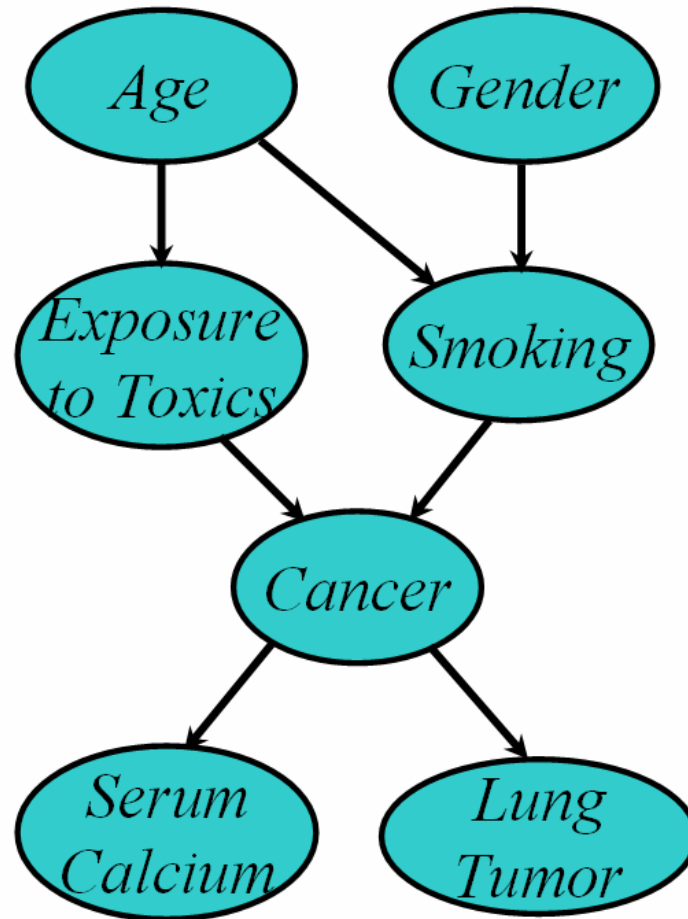
Bayes Rule Revisited

$$P(S|C) = \frac{P(C|S)P(S)}{P(C)} = \frac{P(C,S)}{P(C)}$$

$S \downarrow C \Rightarrow$	<i>none</i>	<i>benign</i>	<i>malig</i>
<i>no</i>	0.768/.935	0.024/.046	0.008/.019
<i>light</i>	0.132/.935	0.012/.046	0.006/.019
<i>heavy</i>	0.030/.935	0.015/.046	0.005/.019

<i>Cancer=</i>	<i>none</i>	<i>benign</i>	<i>malignant</i>
$P(S=no)$	0.821	0.522	0.421
$P(S=light)$	0.141	0.261	0.316
$P(S=heavy)$	0.037	0.217	0.263

A Bayesian Network



Independence

Age

Gender

Age and Gender are independent.

Independence

Age

Gender

Age and Gender are independent.

$$P(A, G) = P(G)P(A)$$

Independence



Age and Gender are independent.

$$P(A, G) = P(G)P(A)$$

$$P(A|G) = P(A) \quad A \perp G$$

$$P(G|A) = P(G) \quad G \perp A$$

Independence

Age

Gender

Age and Gender are independent.

$$P(A, G) = P(G)P(A)$$

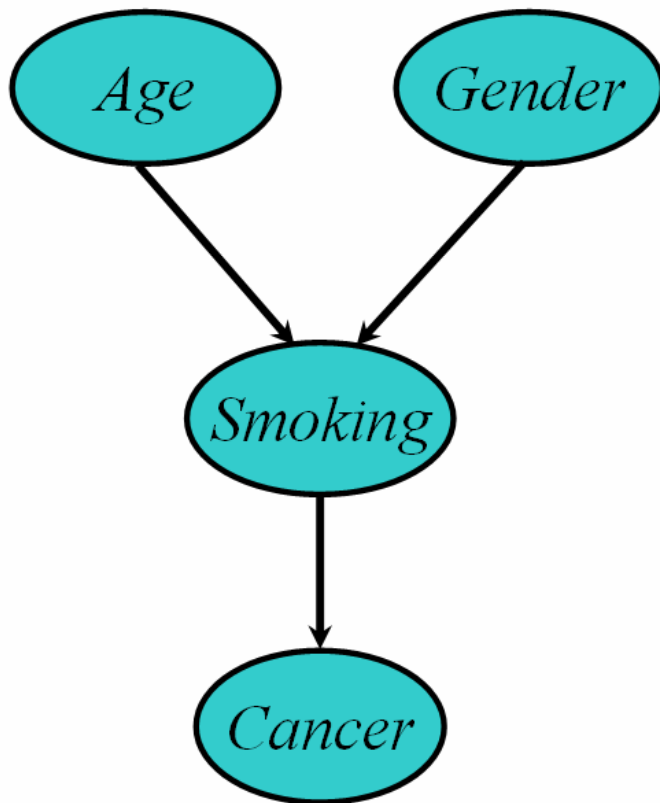
$$P(A|G) = P(A) \quad A \perp G$$

$$P(G|A) = P(G) \quad G \perp A$$

$$P(A, G) = P(G|A) P(A) = P(G)P(A)$$

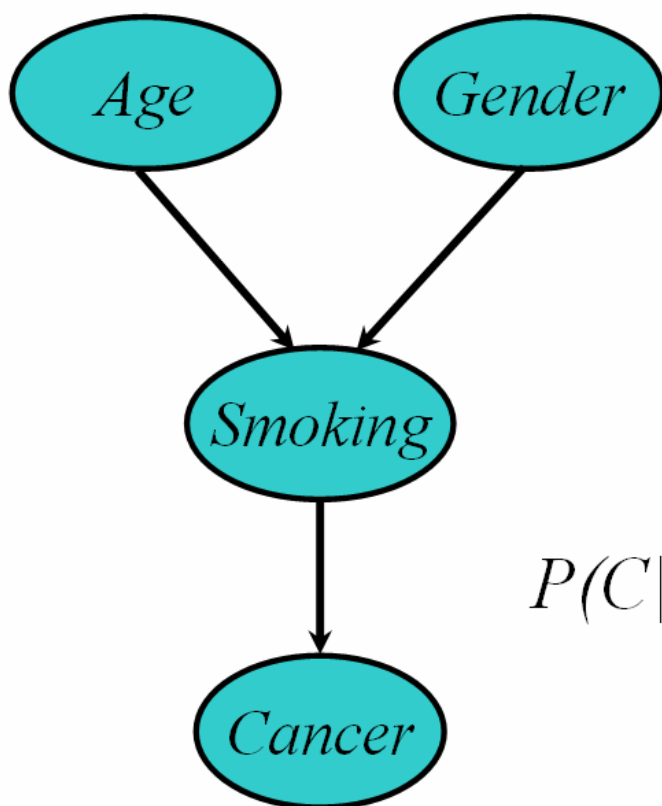
$$P(A, G) = P(A|G) P(G) = P(A)P(G)$$

Conditional Independence



Cancer is independent of Age and Gender given Smoking.

Conditional Independence

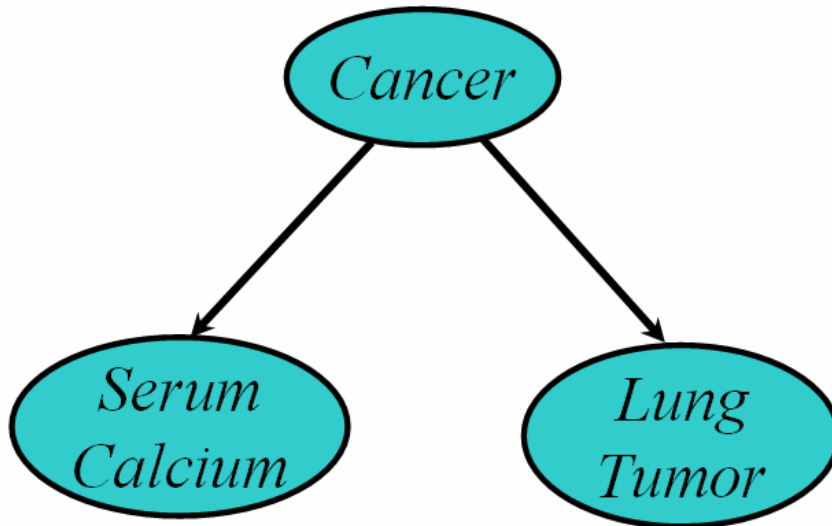


Cancer is independent of *Age* and *Gender* given *Smoking*.

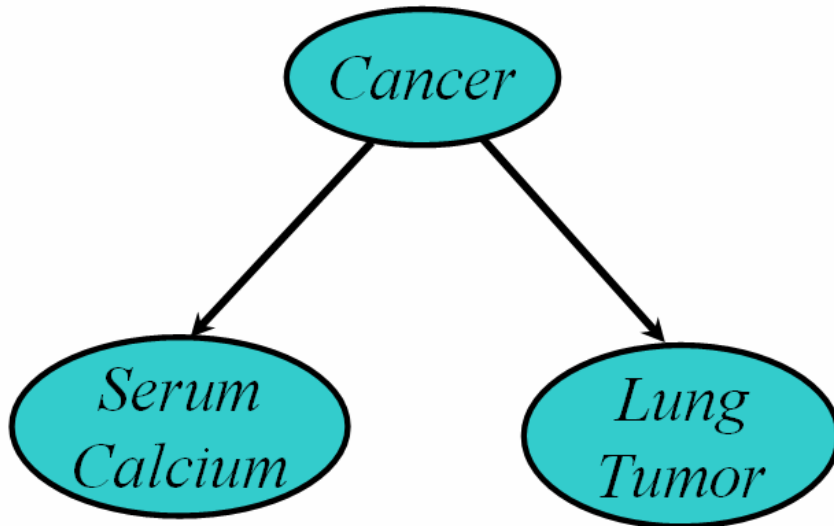
$$P(C|A,G,S) = P(C|S) \quad C \perp A,G \mid S$$

More Conditional Independence: Naïve Bayes

Serum Calcium and Lung Tumor are dependent



More Conditional Independence: Naïve Bayes

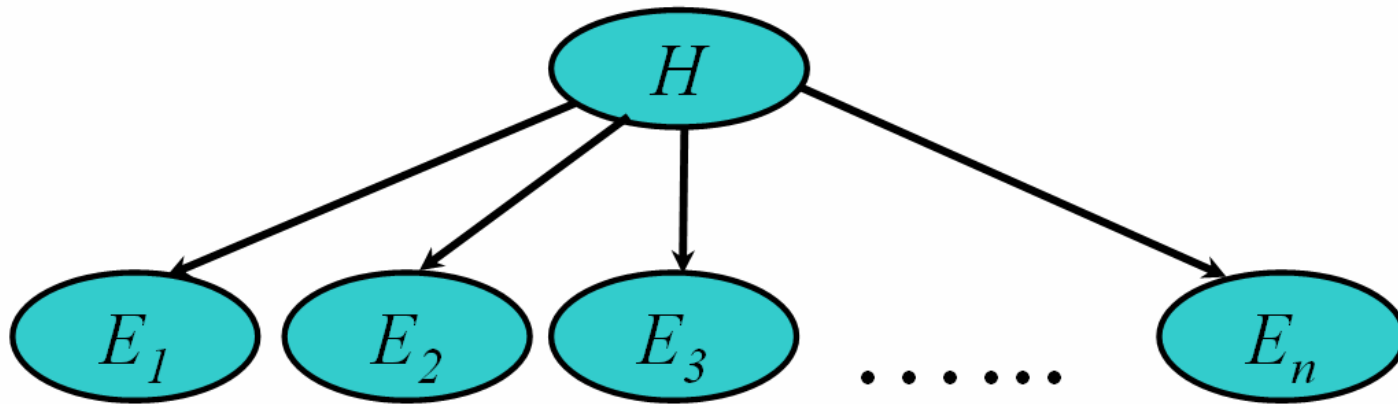


Serum Calcium and Lung Tumor are dependent

Serum Calcium is independent of Lung Tumor, given Cancer

$$P(L|SC, C) = P(L|C)$$

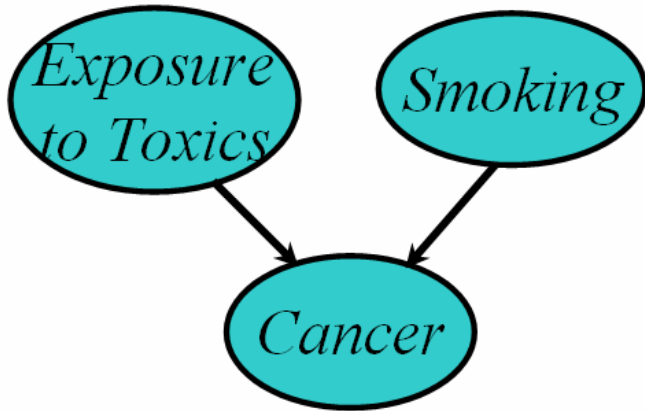
Naïve Bayes in general



$2n + 1$ parameters:

$$P(h)$$
$$P(e_i | h), P(e_i | \bar{h}), \quad i = 1, \dots, n$$

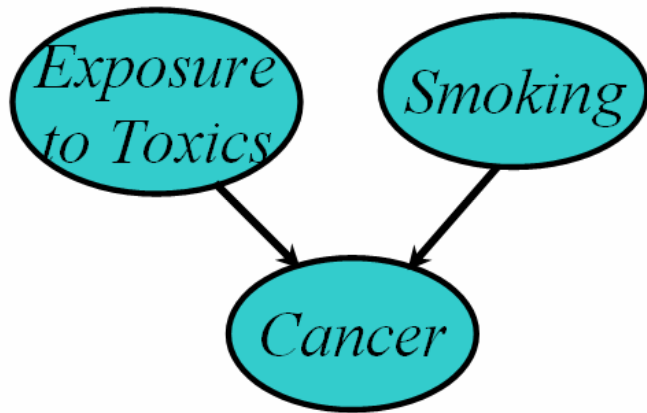
More Conditional Independence: Explaining Away



*Exposure to Toxics and
Smoking are independent*

$$E \perp S$$

More Conditional Independence: Explaining Away

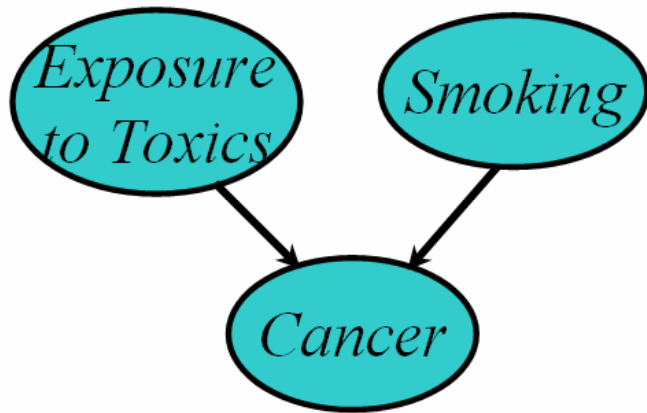


*Exposure to Toxics and
Smoking are independent*

$$E \perp S$$

*Exposure to Toxics is
dependent on *Smoking*,
given *Cancer**

More Conditional Independence: Explaining Away



Exposure to Toxics and Smoking are independent

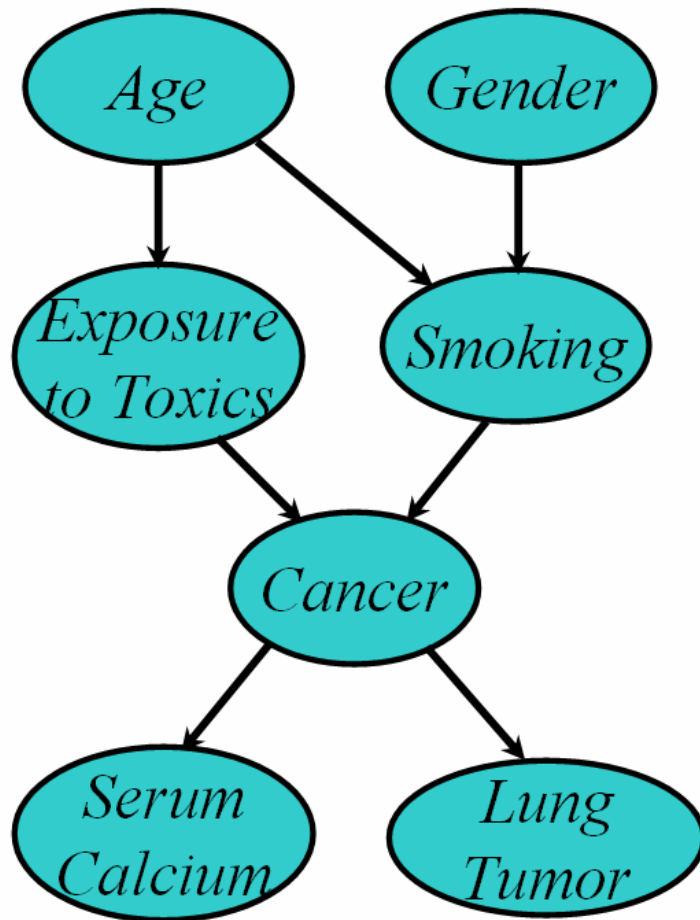
$$E \perp S$$

*Exposure to Toxics is **dependent** on Smoking, given Cancer*

$$P(E = \text{heavy} \mid C = \text{malignant}) >$$

$$P(E = \text{heavy} \mid C = \text{malignant}, S = \text{heavy})$$

Put it all together



$$P(A, G, E, S, C, L, SC) = P(A) \cdot P(G) \cdot$$

$$P(E | A) \cdot P(S | A, G) \cdot$$

$$P(C | E, S) \cdot$$

$$P(SC | C) \cdot P(L | C)$$

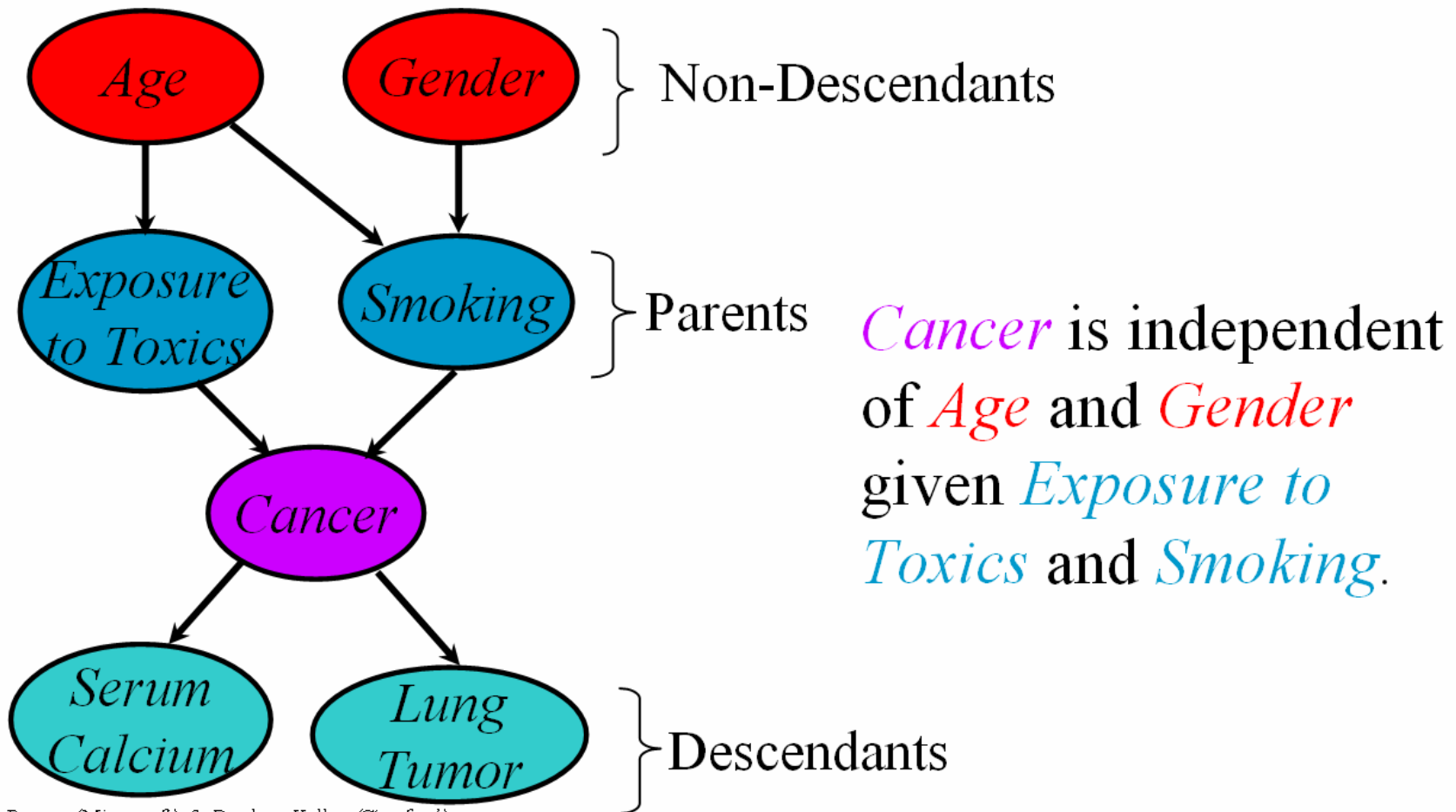
General Product (Chain) Rule for Bayesian Networks

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i \mid \mathbf{Pa}_i)$$

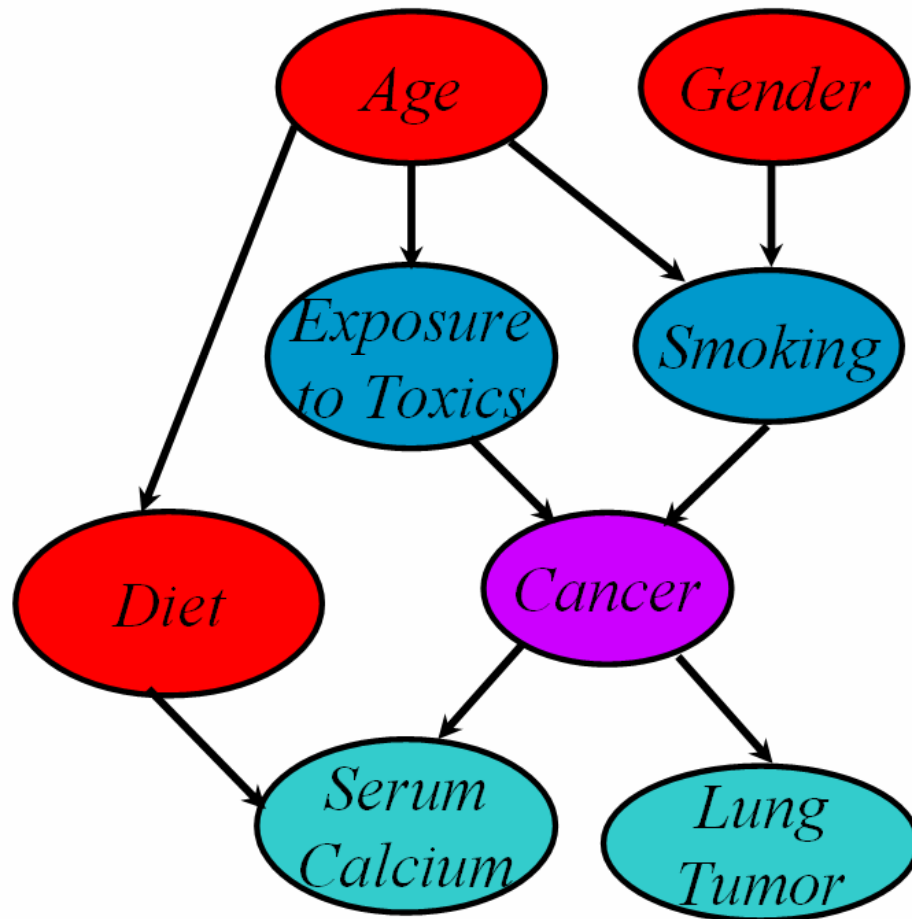
$$\mathbf{Pa}_i = \text{parents}(X_i)$$

Conditional Independence

A variable (node) is conditionally independent of its non-descendants given its parents.



Another non-descendant

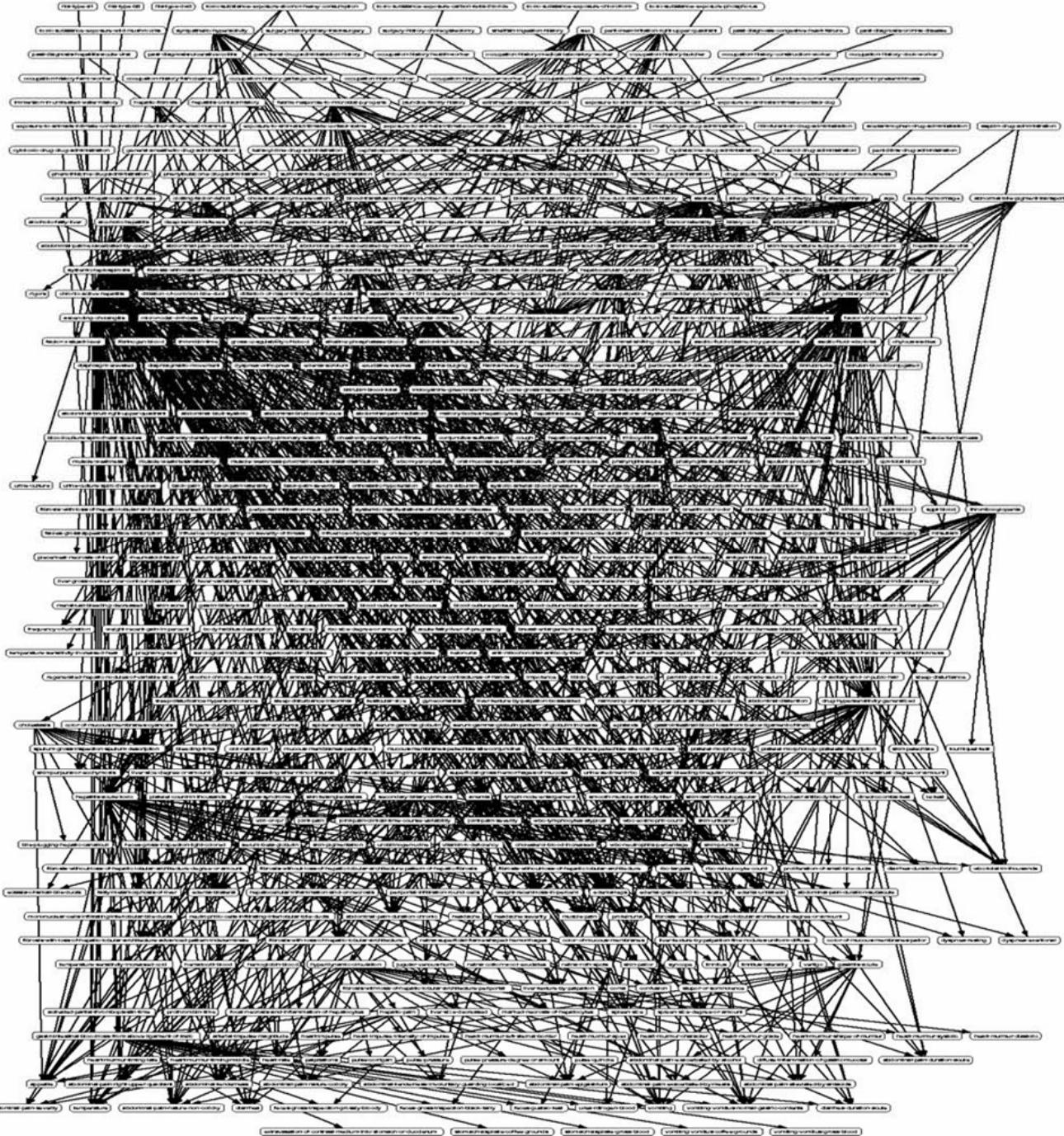


Cancer is independent of *Diet* given *Exposure to Toxics* and *Smoking*.

Independence and Graph Separation

- Given a set of observations, is one set of variables dependent on another set?
- Observing effects can induce dependencies.
- d-separation (Pearl 1988) allows us to check conditional independence graphically.

CPCS Network

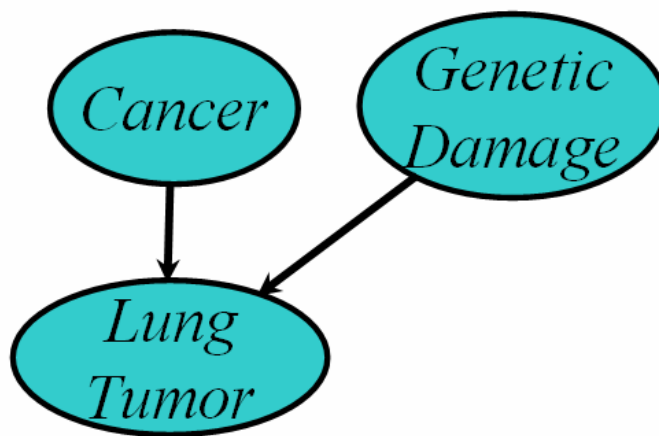


Structuring

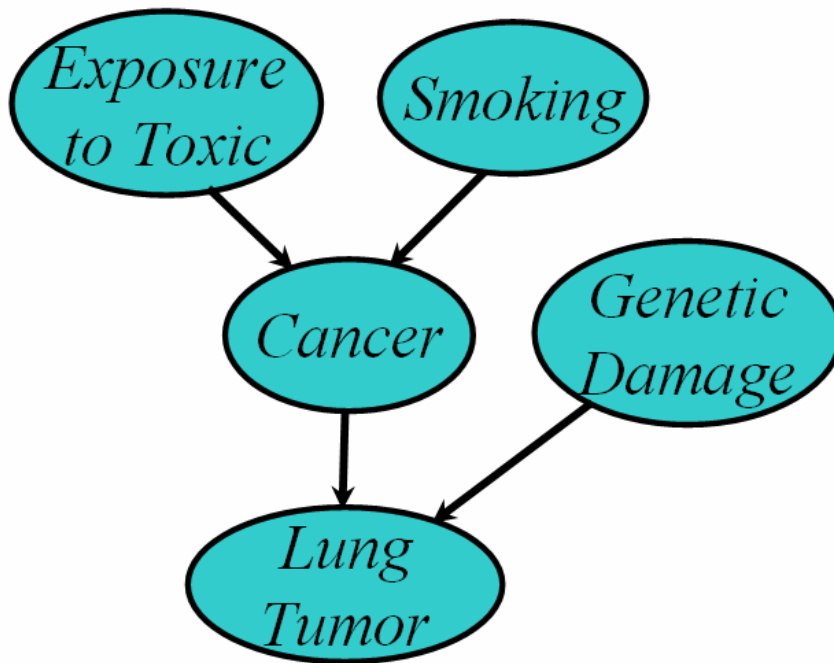
Structuring



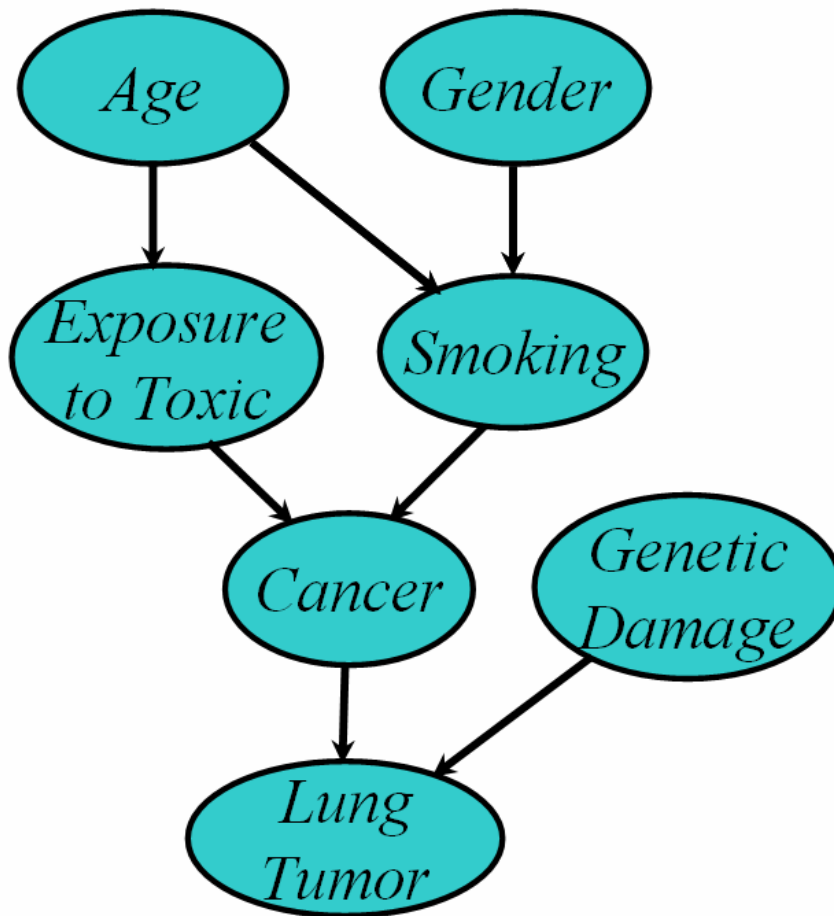
Structuring



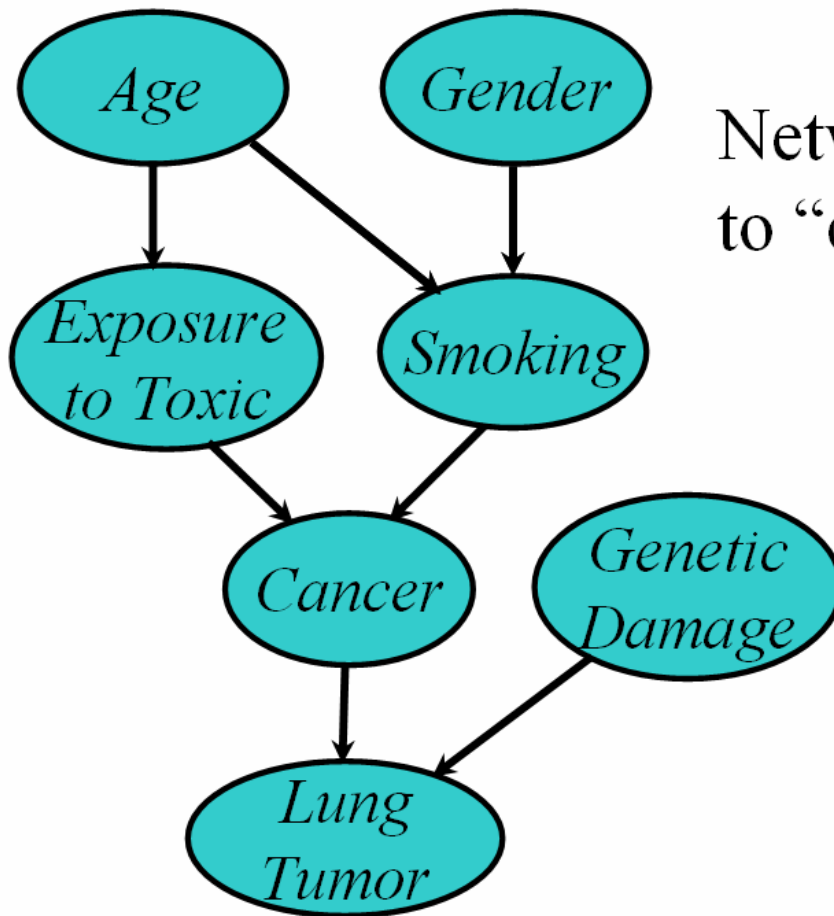
Structuring



Structuring

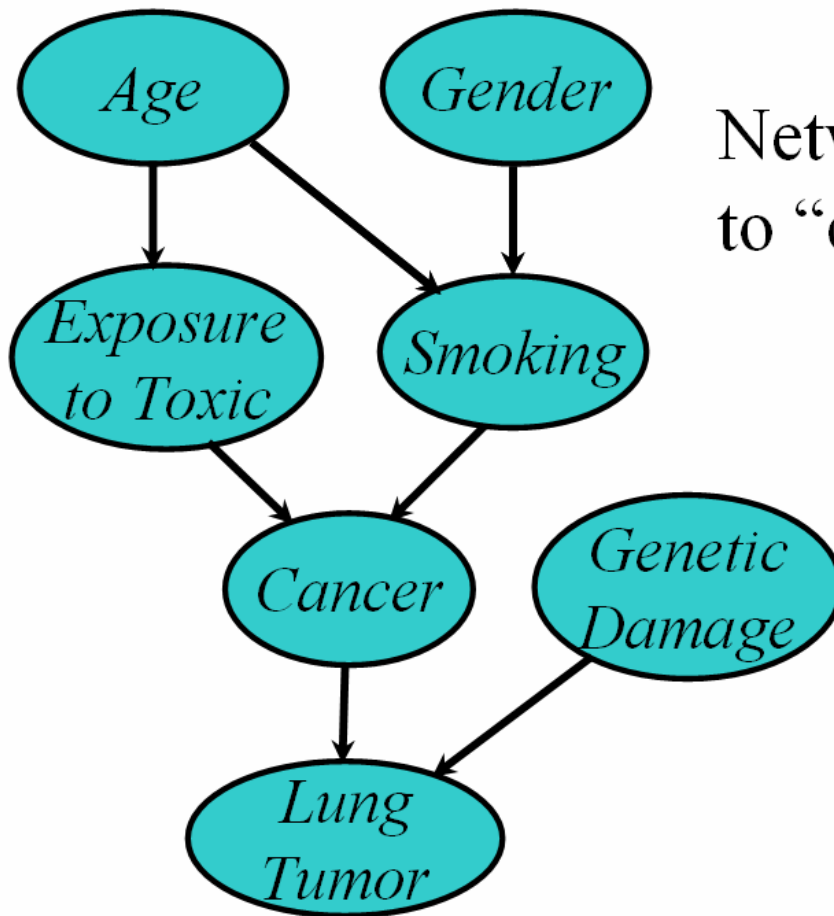


Structuring



Network structure corresponding to “causality” is usually good.

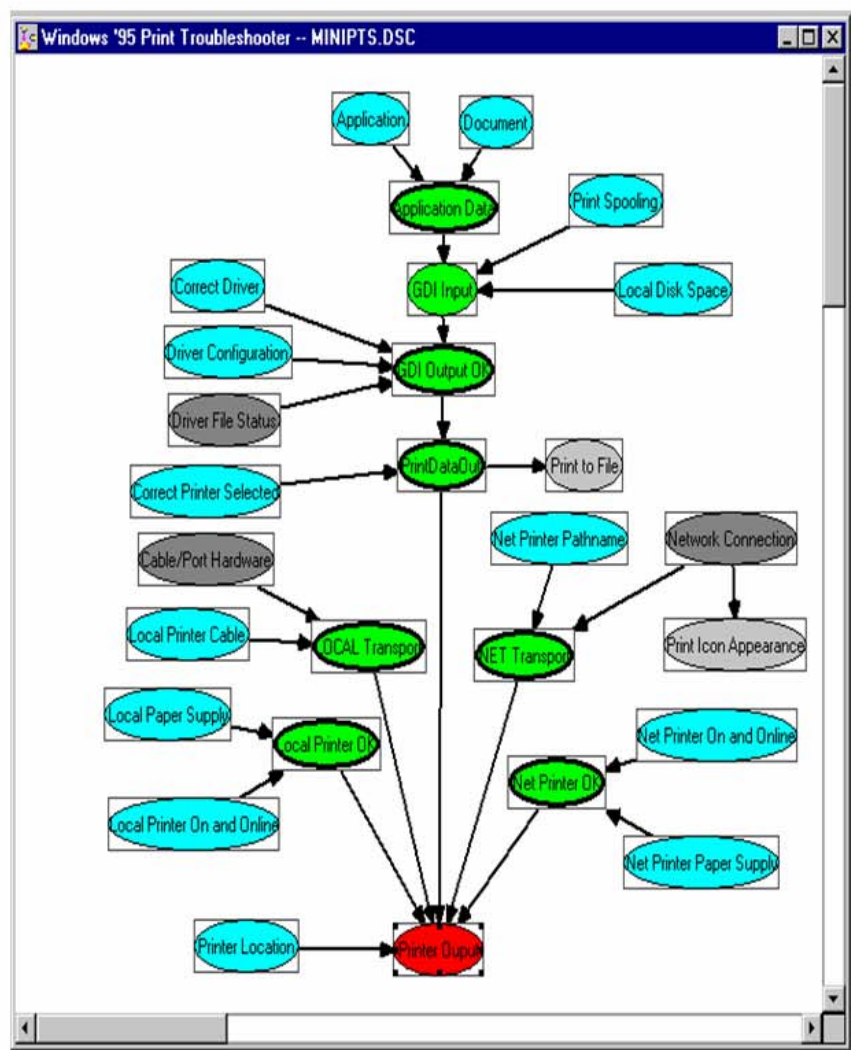
Structuring



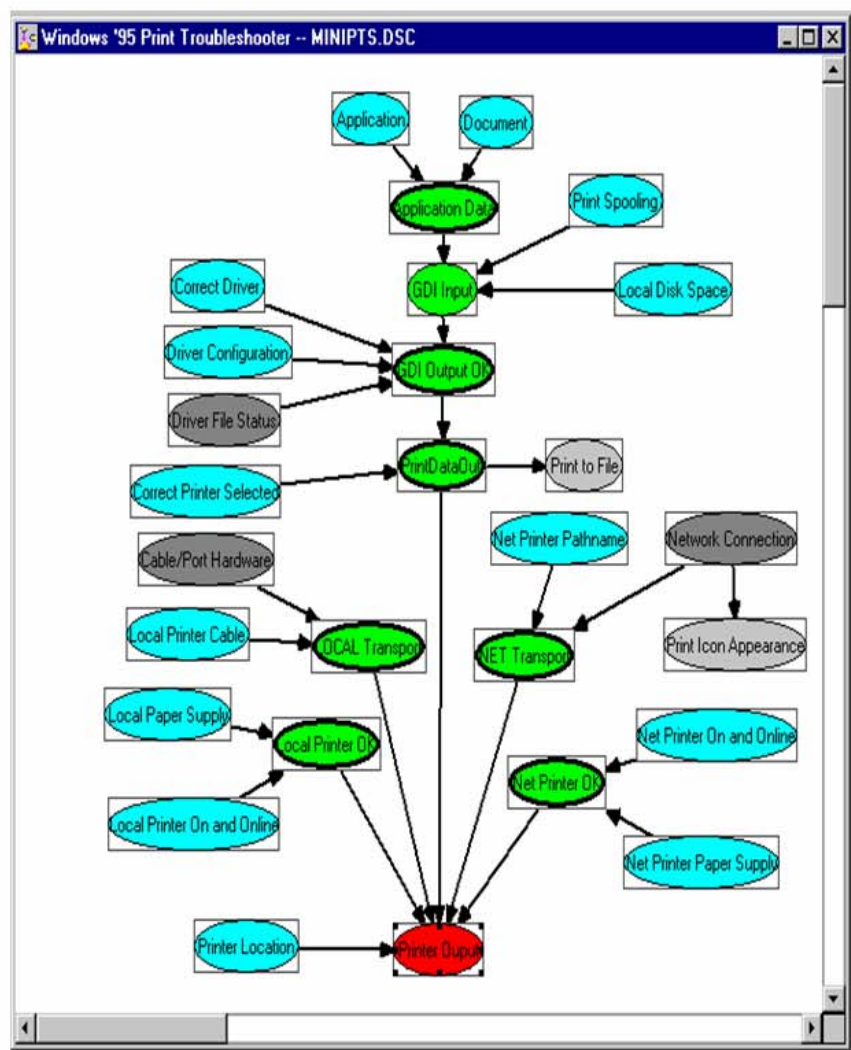
Network structure corresponding to “causality” is usually good.

Extending the conversation.

Local Structure

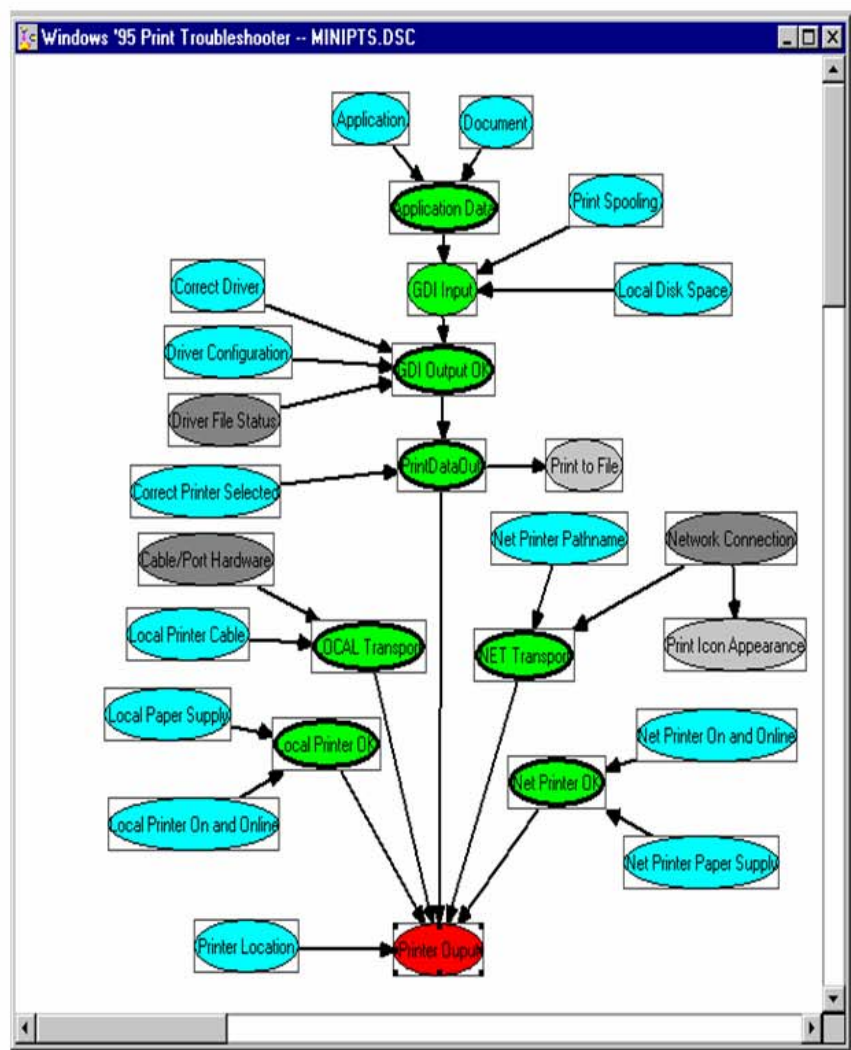


Local Structure



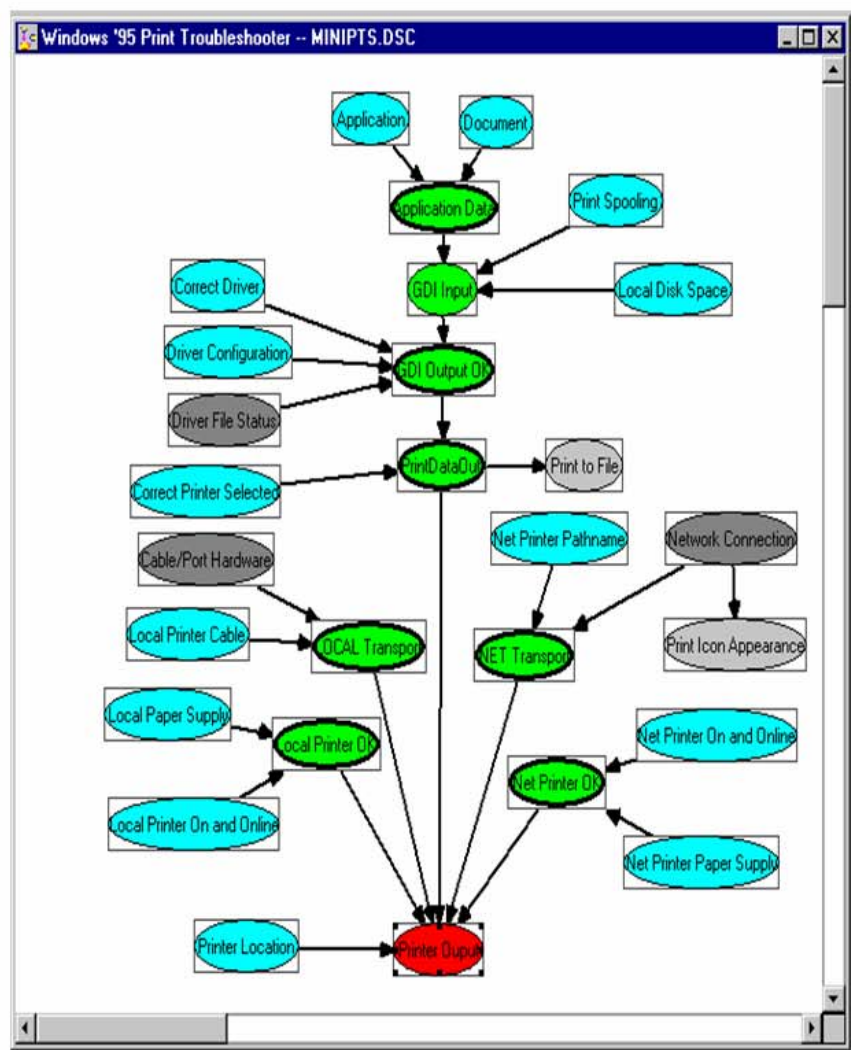
- Causal independence: from 2^n to $n+1$ parameters

Local Structure



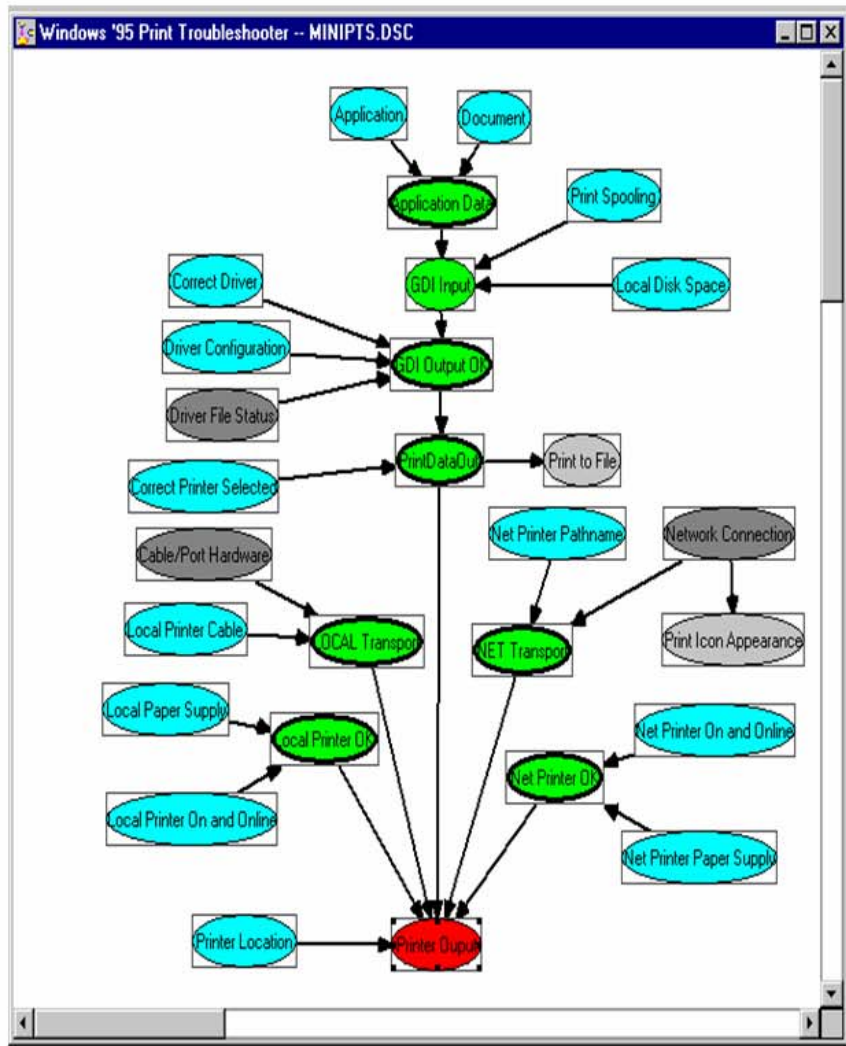
- Causal independence: from 2^n to $n+1$ parameters
- Asymmetric assessment: similar savings in practice.

Local Structure



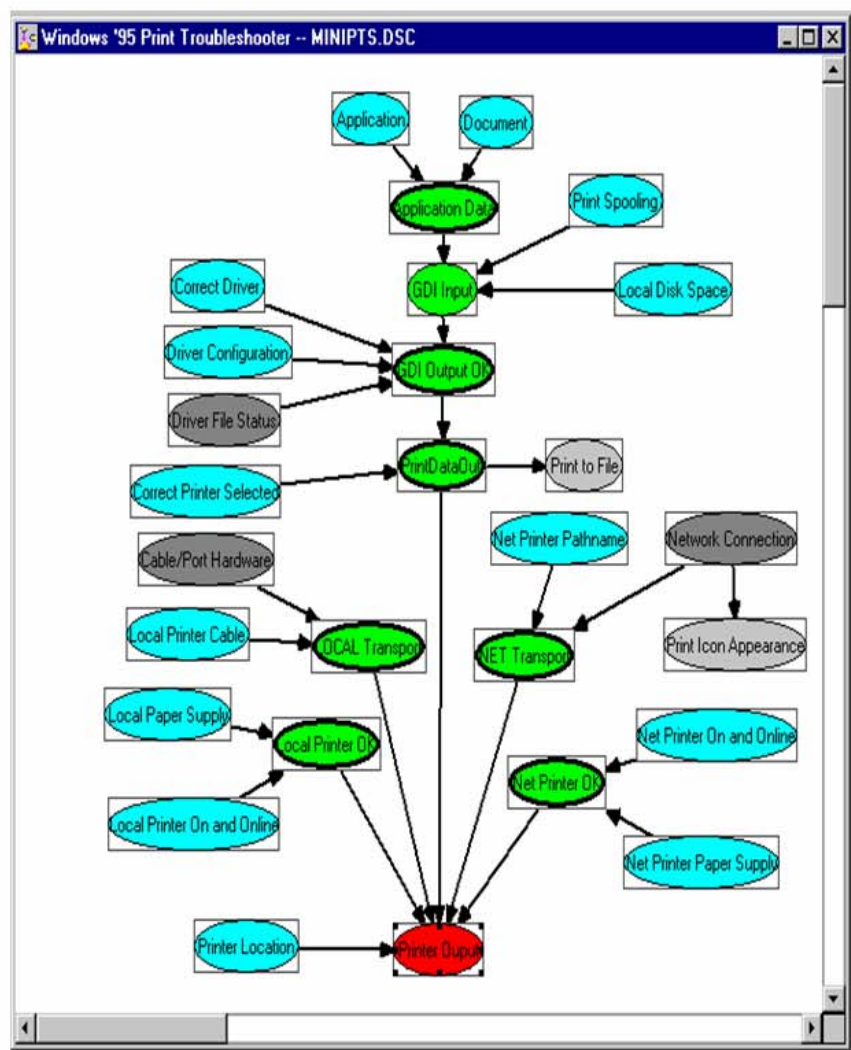
- Causal independence: from 2^n to $n+1$ parameters
- Asymmetric assessment: similar savings in practice.
- Typical savings (#params):

Local Structure



- Causal independence: from 2^n to $n+1$ parameters
- Asymmetric assessment: similar savings in practice.
- Typical savings (#params):
 - ◆ 145 to 55 for a small hardware network;

Local Structure



- Causal independence: from 2^n to $n+1$ parameters
- Asymmetric assessment: similar savings in practice.
- Typical savings (#params):
 - ◆ 145 to 55 for a small hardware network;
 - ◆ 133,931,430 to 8254 for CPCS !!

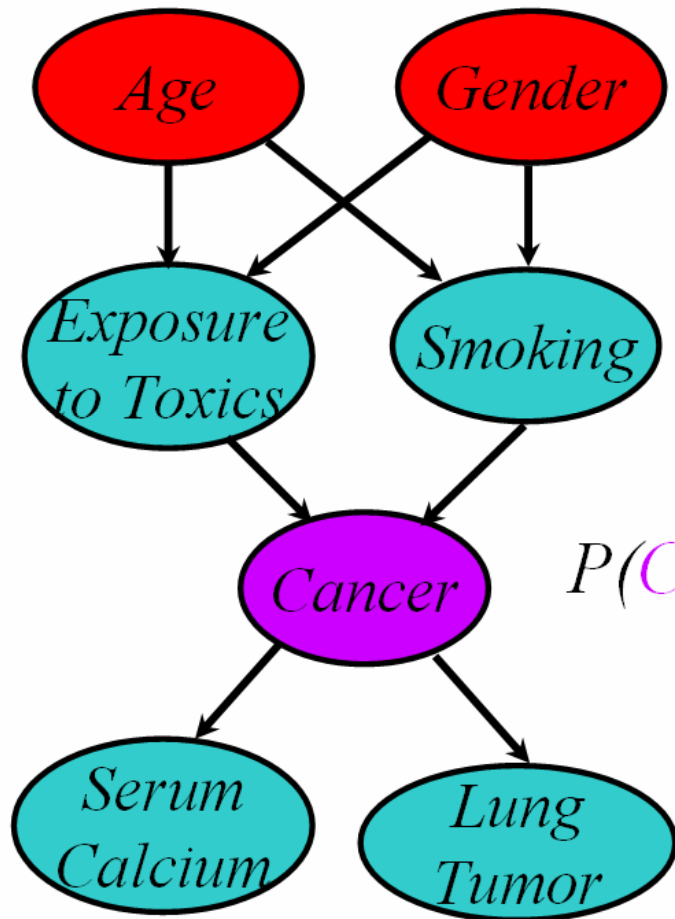
Course Contents

- Concepts in Probability
- Bayesian Networks
- » Inference
- Decision making
- Learning networks from data
- Reasoning over time
- Applications

Inference

- Patterns of reasoning
- Basic inference
- Exact inference
- Exploiting structure
- Approximate inference

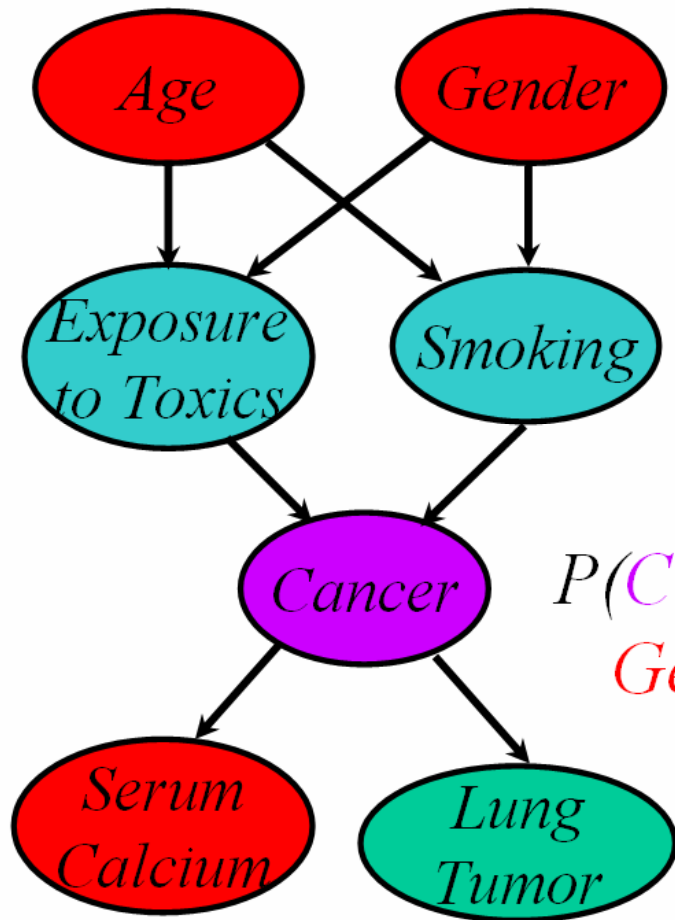
Predictive Inference



How likely are **elderly males** to get **malignant cancer**?

$$P(C=\text{malignant} \mid \text{Age} > 60, \text{Gender} = \text{male})$$

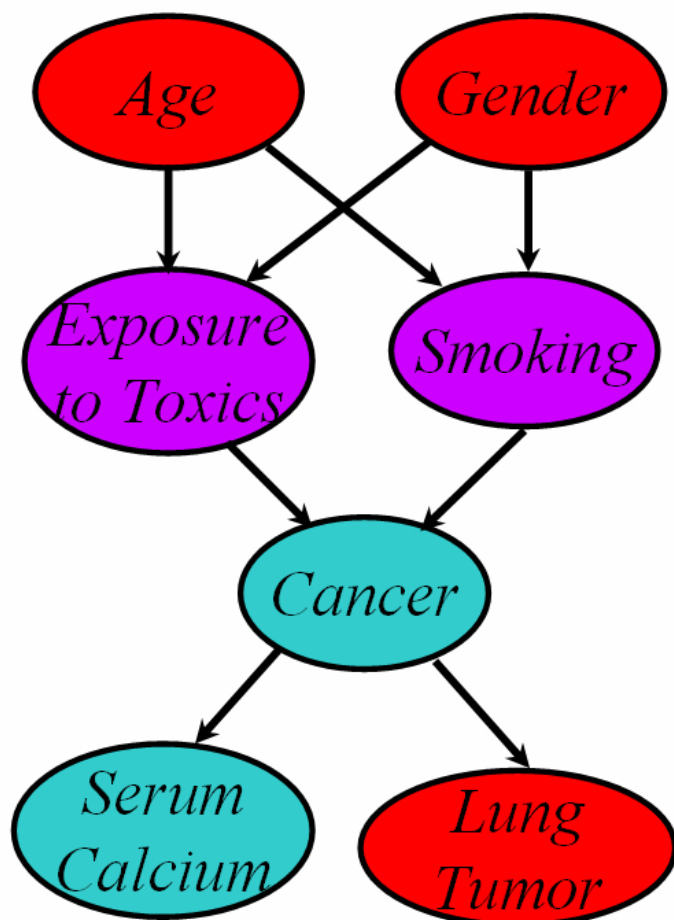
Combined



How likely is an **elderly male** patient with high **Serum Calcium** to have **malignant cancer**?

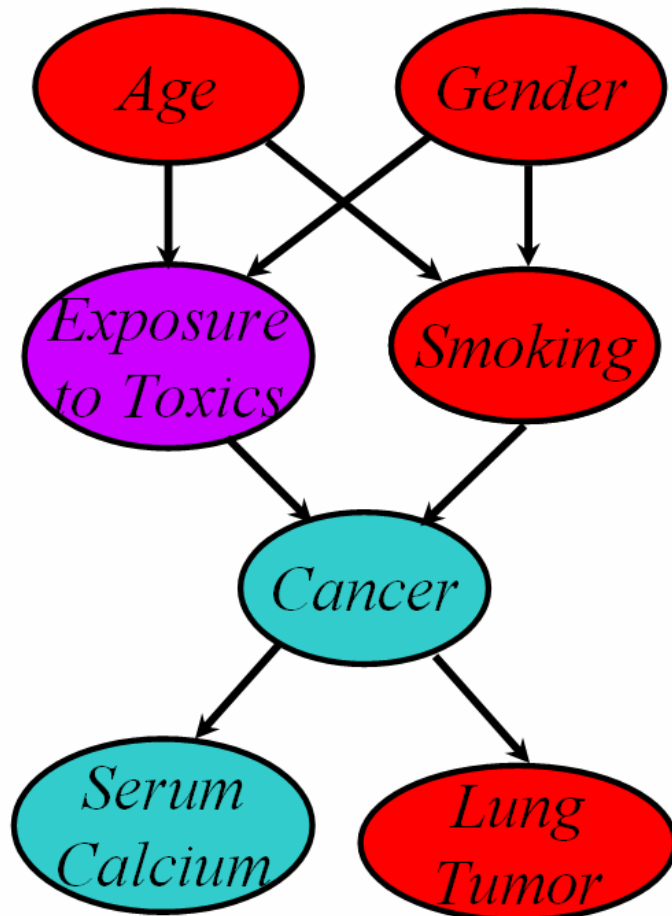
$P(C=\text{malignant} \mid \text{Age} > 60, \text{Gender} = \text{male}, \text{Serum Calcium} = \text{high})$

Explaining away



- If we see a **lung tumor**, the probability of **heavy smoking** and of **exposure to toxics** both go up.

Explaining away



- If we see a **lung tumor**, the probability of **heavy smoking** and of **exposure to toxics** both go up.
- If we then observe **heavy smoking**, the probability of **exposure to toxics** goes back down.

Inference in Belief Networks

- Find $P(Q=q | E=e)$
 - ◆ Q the query variable
 - ◆ E set of evidence variables

$$P(q | e) = \frac{P(q, e)}{P(e)}$$

X_1, \dots, X_n are network variables except Q, E

$$P(q, e) = \sum_{x_1, \dots, x_n} P(q, e, x_1, \dots, x_n)$$

Basic Inference



Basic Inference

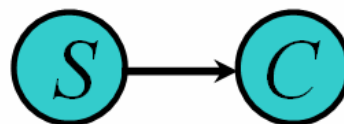
$$P(c) = ?$$


```
graph LR; S((S)) --> C((C))
```

- $P(C, S) = P(C|S) P(S)$

Basic Inference

$$P(c) = ?$$



- $P(C,S) = P(C|S) P(S)$

<i>Smoking=</i>	<i>no</i>	<i>light</i>	<i>heavy</i>
$P(C=none)$	0.96	0.88	0.60
$P(C=benign)$	0.03	0.08	0.25
$P(C=malig)$	0.01	0.04	0.15

Basic Inference



■ $P(C,S) = P(C|S) P(S)$

<i>Smoking=</i>	<i>no</i>	<i>light</i>	<i>heavy</i>
$P(C=none)$	0.96	0.88	0.60
$P(C=benign)$	0.03	0.08	0.25
$P(C=malig)$	0.01	0.04	0.15

$P(S=no)$	0.80
$P(S=light)$	0.15
$P(S=heavy)$	0.05

Basic Inference



■ $P(C,S) = P(C|S) P(S)$

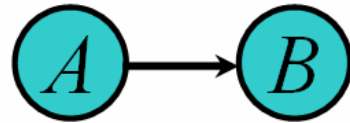
<i>Smoking=</i>	<i>no</i>	<i>light</i>	<i>heavy</i>
$P(C=none)$	0.96	0.88	0.60
$P(C=benign)$	0.03	0.08	0.25
$P(C=malig)$	0.01	0.04	0.15

$P(S=no)$	0.80
$P(S=light)$	0.15
$P(S=heavy)$	0.05

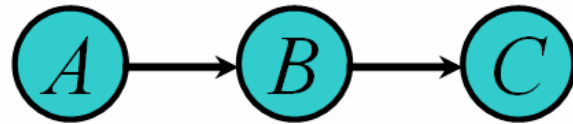
$S \Downarrow$ $C \Rightarrow$	<i>none</i>	<i>benign</i>	<i>malig</i>	total
<i>no</i>	0.768	0.024	0.008	.80
<i>light</i>	0.132	0.012	0.006	.15
<i>heavy</i>	0.035	0.010	0.005	.05
total	0.935	0.046	0.019	


 $P(\text{Cancer})$

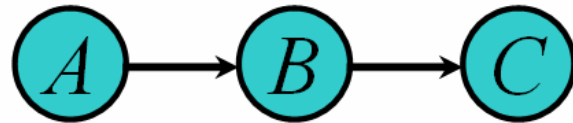
Basic Inference



Basic Inference

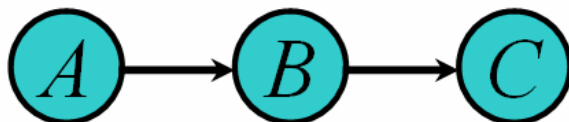


Basic Inference



$$P(b) = \sum_a P(a, b) = \sum_a P(b | a) P(a)$$

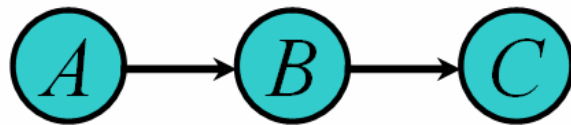
Basic Inference



$$P(b) = \sum_a P(a, b) = \sum_a P(b | a) P(a)$$

$$P(c) = \sum_b P(c | b) P(b)$$

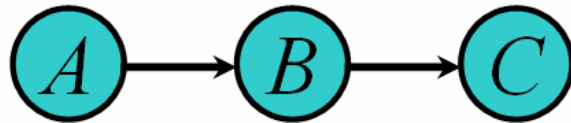
Basic Inference



$$\underbrace{P(b)} = \sum_a P(a, b) = \sum_a P(b | a) P(a)$$

$$P(c) = \sum_b P(c | b) \overbrace{P(b)}$$

Basic Inference

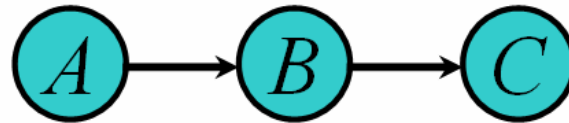


$$\underbrace{P(b)} = \sum_a P(a, b) = \sum_a P(b | a) P(a)$$

$$P(c) = \sum_b P(c | b) \overbrace{P(b)}$$

$$P(c) = \sum_{b,a} P(a, b, c)$$

Basic Inference

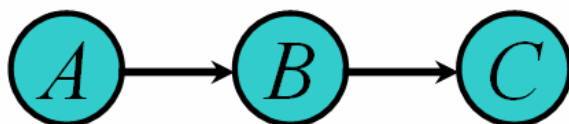


$$\underbrace{P(b)} = \sum_a P(a, b) = \sum_a P(b | a) P(a)$$

$$P(c) = \sum_b P(c | b) \overbrace{P(b)}$$

$$P(c) = \sum_{b,a} P(a, b, c) = \sum_{b,a} P(c | b) P(b | a) P(a)$$

Basic Inference



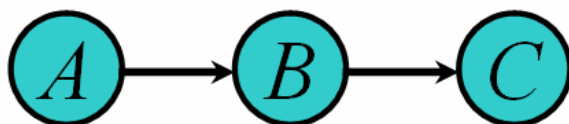
$$\underbrace{P(b)} = \sum_a P(a, b) = \sum_a P(b | a) P(a)$$

$$P(c) = \sum_b P(c | b) \underbrace{P(b)}$$

$$P(c) = \sum_{b,a} P(a, b, c) = \sum_{b,a} P(c | b) P(b | a) P(a)$$

$$= \sum_b P(c | b) \sum_a P(b | a) P(a)$$

Basic Inference



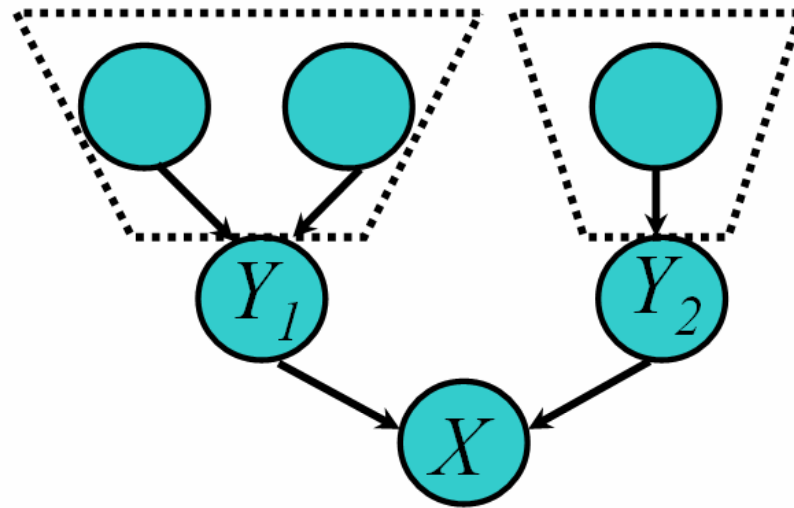
$$\underbrace{P(b)} = \sum_a P(a, b) = \sum_a P(b | a) P(a)$$

$$P(c) = \sum_b P(c | b) \underbrace{P(b)}$$

$$P(c) = \sum_{b,a} P(a, b, c) = \sum_{b,a} P(c | b) P(b | a) P(a)$$

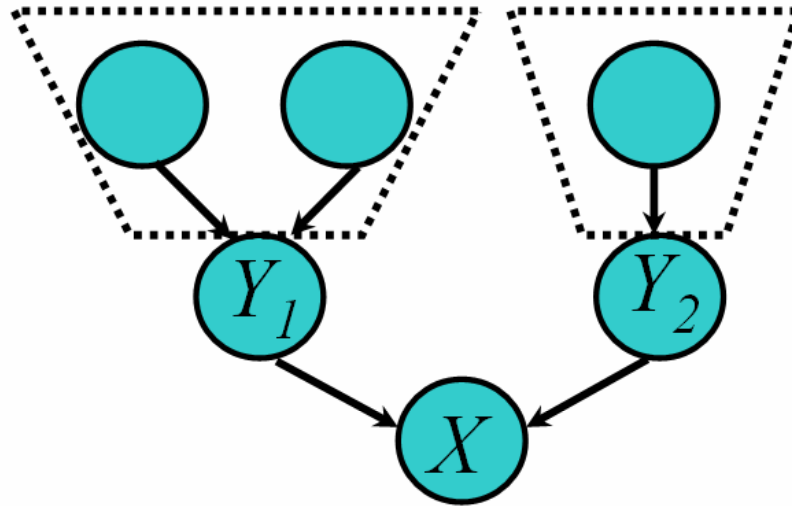
$$= \sum_b P(c | b) \underbrace{\sum_a P(b | a) P(a)}_{P(b)}$$

Inference in trees



$$P(x) = \sum_{y_1, y_2} P(x | y_1, y_2) P(y_1, y_2)$$

Inference in trees



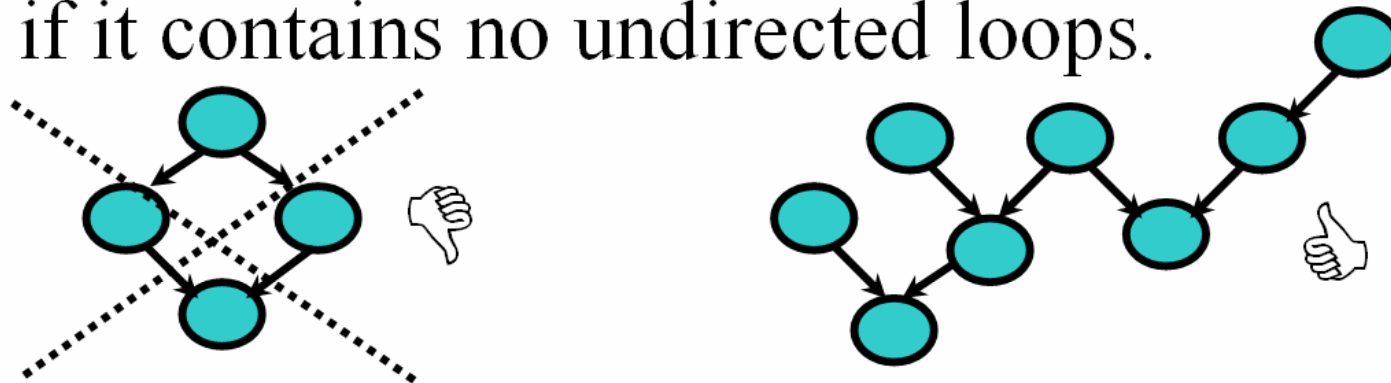
$$P(x) = \sum_{y_1, y_2} P(x | y_1, y_2) P(y_1, y_2)$$

because of independence of Y_1, Y_2 :

$$= \sum_{y_1, y_2} P(x | y_1, y_2) P(y_1) P(y_2)$$

Polytrees

- A network is *singly connected* (a *polytree*) if it contains no undirected loops.

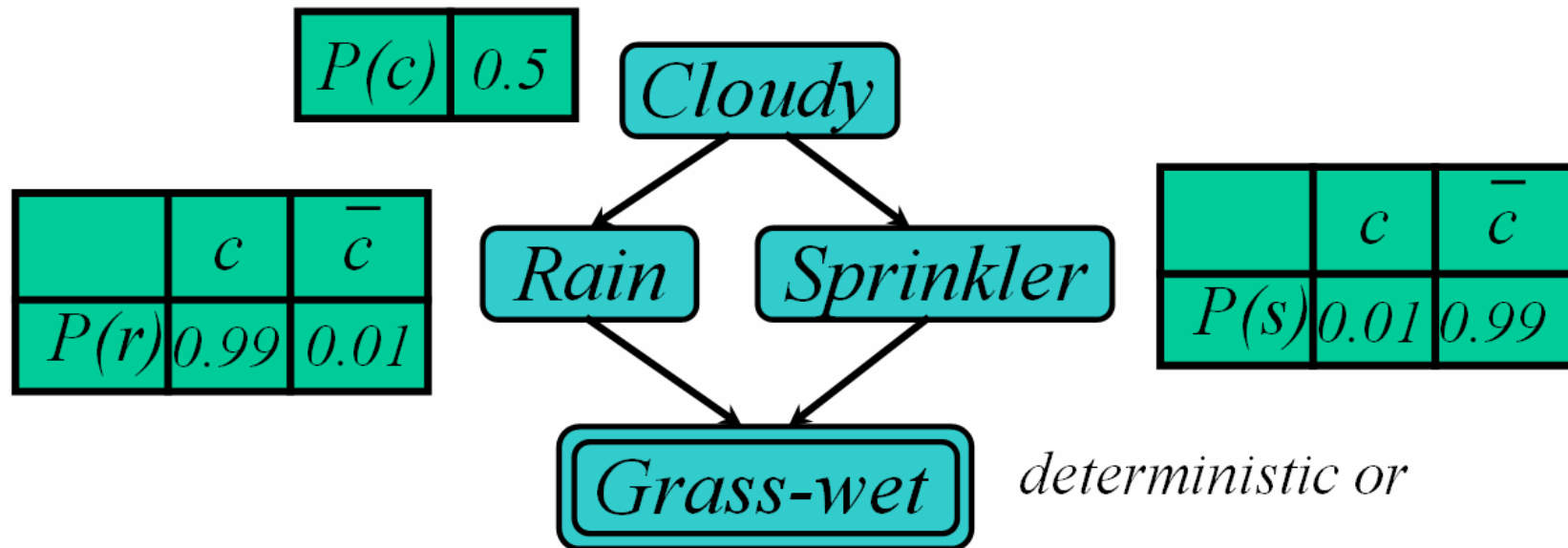


Theorem: Inference in a singly connected network can be done in linear time*.

Main idea: in variable elimination, need only maintain distributions over single nodes.

* in network size including table sizes.

The problem with loops



The grass is dry only if no rain and no sprinklers.

$$P(\bar{g}) = P(\bar{r}, \bar{s}) \sim 0$$

The problem with loops contd.

$$P(\bar{g}) =$$

The problem with loops contd.

$$P(\bar{g}) = P(\bar{g} | r, s) P(r, s) + P(\bar{g} | r, \bar{s}) P(r, \bar{s}) \\ + P(\bar{g} | \bar{r}, s) P(\bar{r}, s) + P(\bar{g} | \bar{r}, \bar{s}) P(\bar{r}, \bar{s})$$

The problem with loops contd.

$$P(\bar{g}) = \underbrace{P(\bar{g} | r, s)}_0 P(r, s) + \underbrace{P(\bar{g} | r, \bar{s})}_0 P(r, \bar{s}) \\ + \underbrace{P(\bar{g} | \bar{r}, s)}_0 P(\bar{r}, s) + \underbrace{P(\bar{g} | \bar{r}, \bar{s})}_1 P(\bar{r}, \bar{s})$$

The problem with loops contd.

$$\begin{aligned} P(\bar{g}) &= \overbrace{P(\bar{g} \mid r, s)}^0 P(r, s) + \overbrace{P(\bar{g} \mid r, \bar{s})}^0 P(r, \bar{s}) \\ &\quad + \underbrace{P(\bar{g} \mid \bar{r}, s)}_0 P(\bar{r}, s) + \underbrace{P(\bar{g} \mid \bar{r}, \bar{s})}_1 P(\bar{r}, \bar{s}) \\ &= P(\bar{r}, \bar{s}) \end{aligned}$$

The problem with loops contd.

$$\begin{aligned} P(\bar{g}) &= \overbrace{P(\bar{g} | r, s)}^0 P(r, s) + \overbrace{P(\bar{g} | r, \bar{s})}^0 P(r, \bar{s}) \\ &\quad + \underbrace{P(\bar{g} | \bar{r}, s)}_0 P(\bar{r}, s) + \underbrace{P(\bar{g} | \bar{r}, \bar{s})}_1 P(\bar{r}, \bar{s}) \\ &= P(\bar{r}, \bar{s}) \\ &= P(\bar{r}) P(\bar{s}) \sim 0.5 \cdot 0.5 = 0.25 \end{aligned}$$

The problem with loops contd.

$$\begin{aligned} P(\bar{g}) &= \overbrace{P(\bar{g} | r, s)}^0 P(r, s) + \overbrace{P(\bar{g} | r, \bar{s})}^0 P(r, \bar{s}) \\ &\quad + \underbrace{P(\bar{g} | \bar{r}, s)}_0 P(\bar{r}, s) + \underbrace{P(\bar{g} | \bar{r}, \bar{s})}_1 P(\bar{r}, \bar{s}) \\ &= P(\bar{r}, \bar{s}) \sim 0 \\ &= P(\bar{r}) P(\bar{s}) \sim 0.5 \cdot 0.5 = 0.25 \end{aligned}$$

The problem with loops contd.

$$P(\bar{g}) = \overbrace{P(\bar{g} | r, s)}^0 P(r, s) + \overbrace{P(\bar{g} | r, \bar{s})}^0 P(r, \bar{s}) \\ + \underbrace{P(\bar{g} | \bar{r}, s)}_0 P(\bar{r}, s) + \underbrace{P(\bar{g} | \bar{r}, \bar{s})}_1 P(\bar{r}, \bar{s})$$

$$= P(\bar{r}, \bar{s}) \sim 0$$

$$\neq P(\bar{r}) P(\bar{s}) \sim 0.5 \cdot 0.5 = 0.25$$

problem

Variable elimination

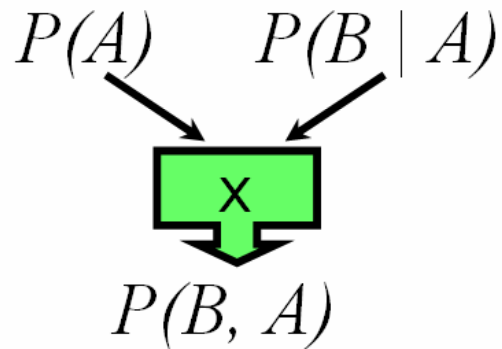


$$P(c) = \sum_b P(c | b) \underbrace{\sum_a P(b | a) P(a)}_{P(b)}$$

Variable elimination



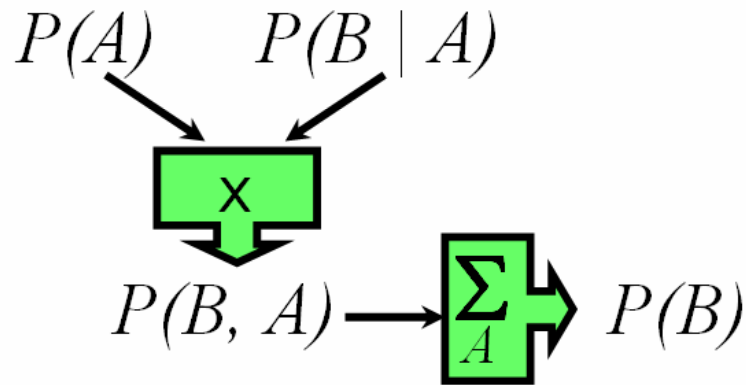
$$P(c) = \sum_b P(c | b) \underbrace{\sum_a P(b | a) P(a)}_{P(b)}$$



Variable elimination



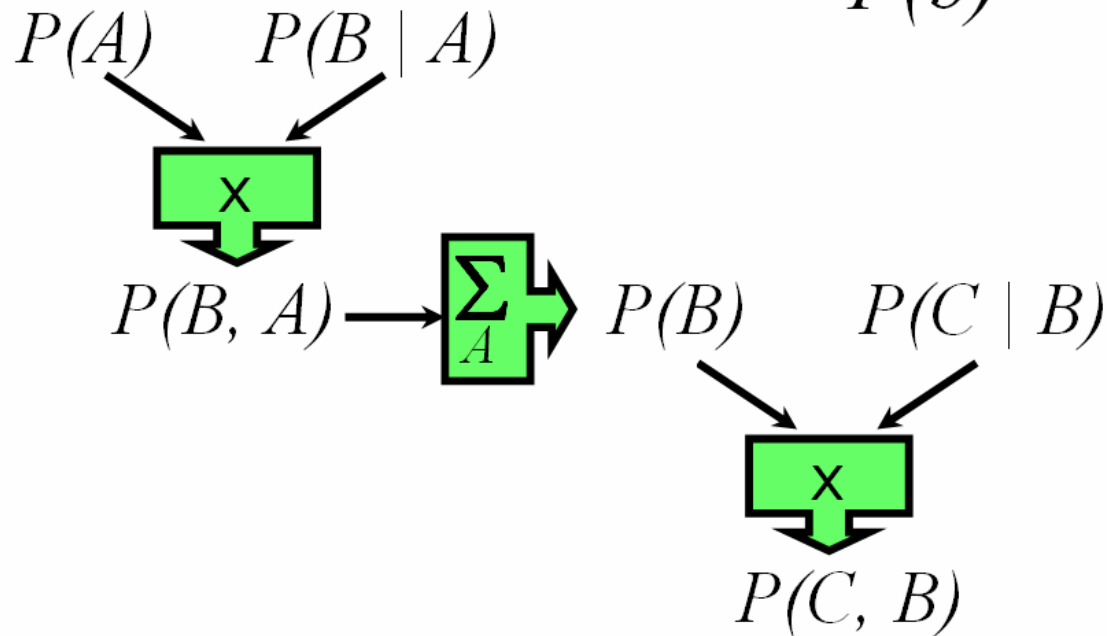
$$P(c) = \sum_b P(c | b) \underbrace{\sum_a P(b | a) P(a)}_{P(b)}$$



Variable elimination



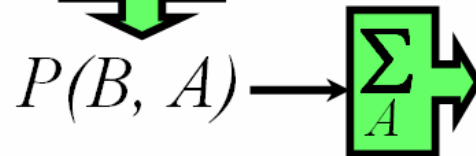
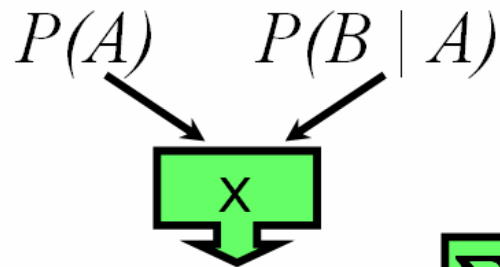
$$P(c) = \sum_b P(c | b) \underbrace{\sum_a P(b | a) P(a)}_{P(b)}$$



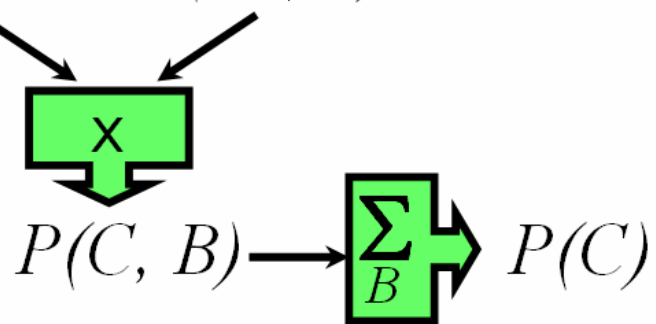
Variable elimination



$$P(c) = \sum_b P(c | b) \underbrace{\sum_a P(b | a) P(a)}_{P(b)}$$



$$P(B) \quad P(C | B)$$



Inference as variable elimination

Inference as variable elimination

- A **factor** over X is a function from $val(X)$ to numbers in $[0, 1]$:

Inference as variable elimination

- A **factor** over X is a function from $val(X)$ to numbers in $[0,1]$:
 - ◆ A CPT is a factor

Inference as variable elimination

- A **factor** over X is a function from $val(X)$ to numbers in $[0, 1]$:
 - ◆ A CPT is a factor
 - ◆ A joint distribution is also a factor

Inference as variable elimination

- A **factor** over X is a function from $val(X)$ to numbers in $[0, 1]$:
 - ◆ A CPT is a factor
 - ◆ A joint distribution is also a factor
- BN inference:

Inference as variable elimination

- A **factor** over X is a function from $val(X)$ to numbers in $[0, 1]$:
 - ◆ A CPT is a factor
 - ◆ A joint distribution is also a factor
- BN inference:
 - ◆ factors are multiplied to give new ones

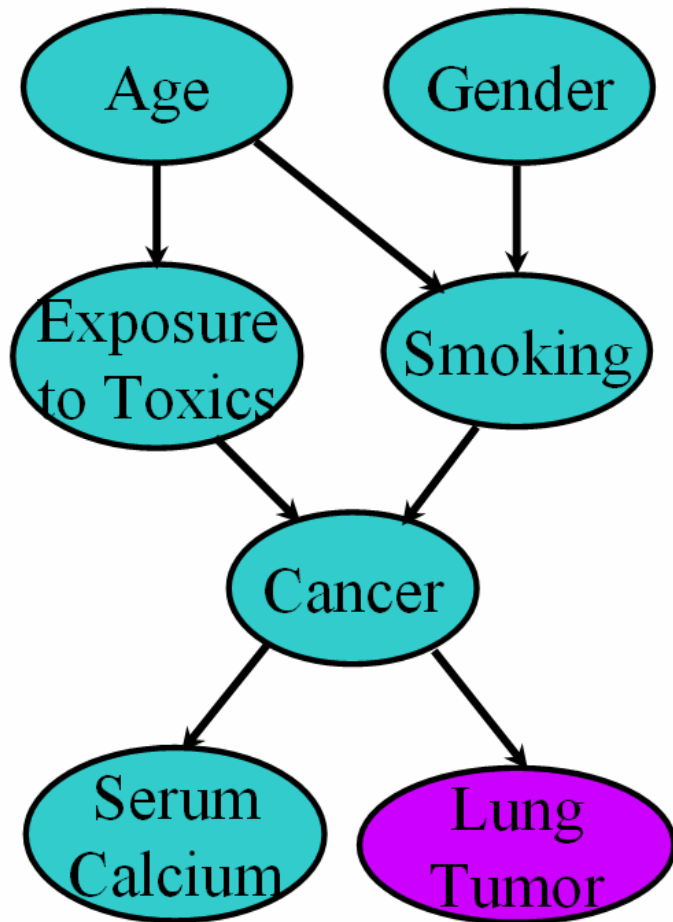
Inference as variable elimination

- A **factor** over X is a function from $val(X)$ to numbers in $[0, 1]$:
 - ◆ A CPT is a factor
 - ◆ A joint distribution is also a factor
- BN inference:
 - ◆ factors are multiplied to give new ones
 - ◆ variables in factors summed out

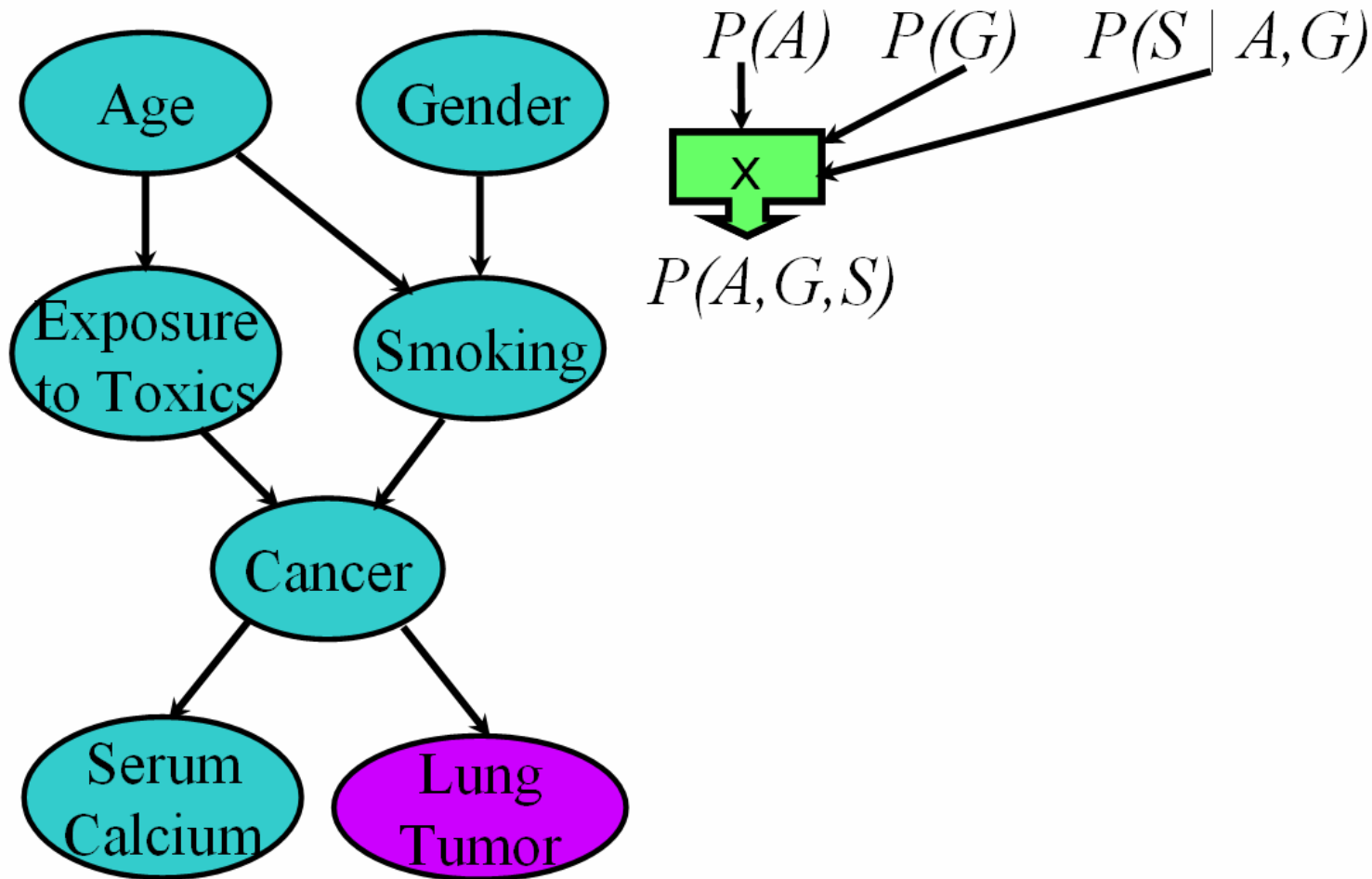
Inference as variable elimination

- A **factor** over X is a function from $val(X)$ to numbers in $[0, 1]$:
 - ◆ A CPT is a factor
 - ◆ A joint distribution is also a factor
- BN inference:
 - ◆ factors are multiplied to give new ones
 - ◆ variables in factors summed out
- A variable can be summed out as soon as all factors mentioning it have been multiplied.

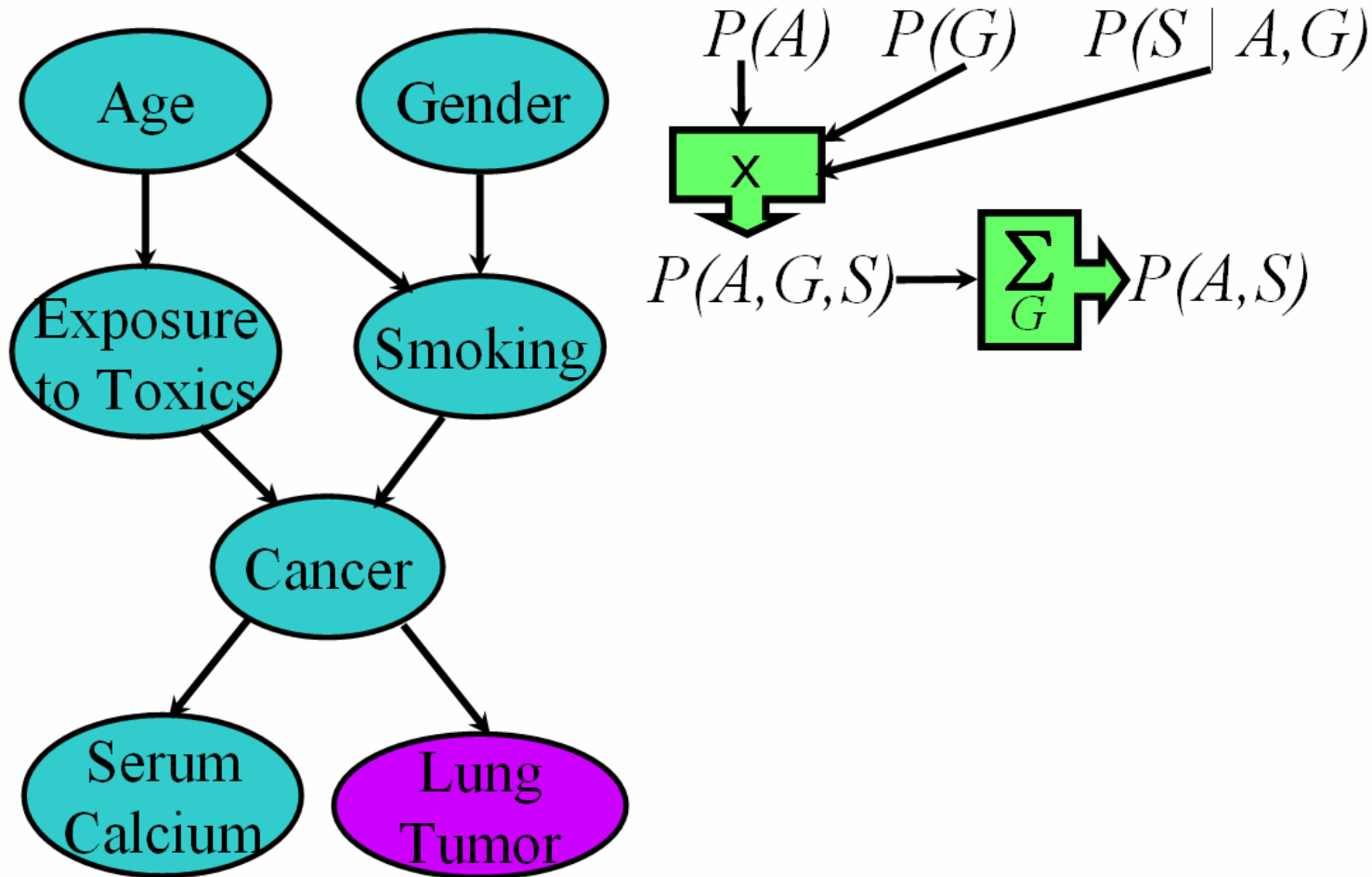
Variable Elimination with loops



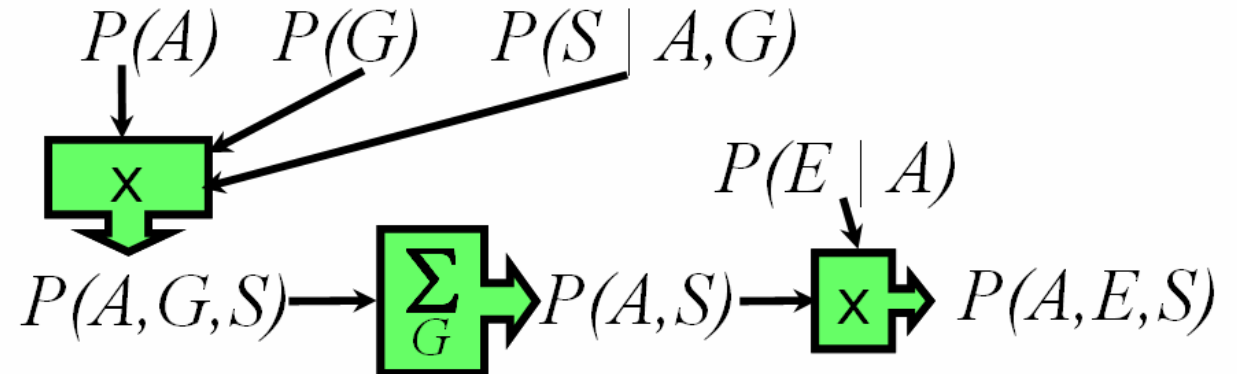
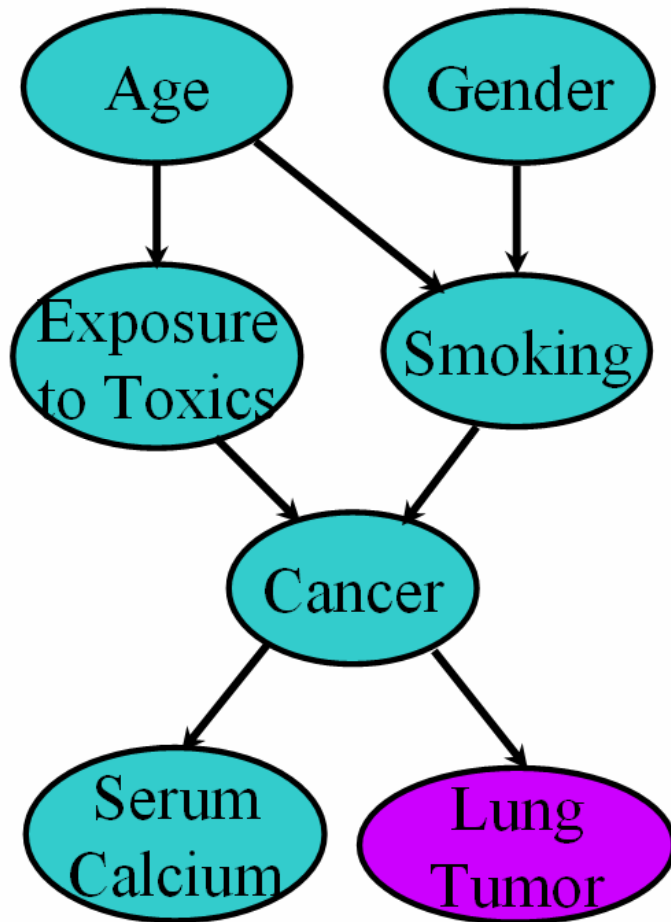
Variable Elimination with loops



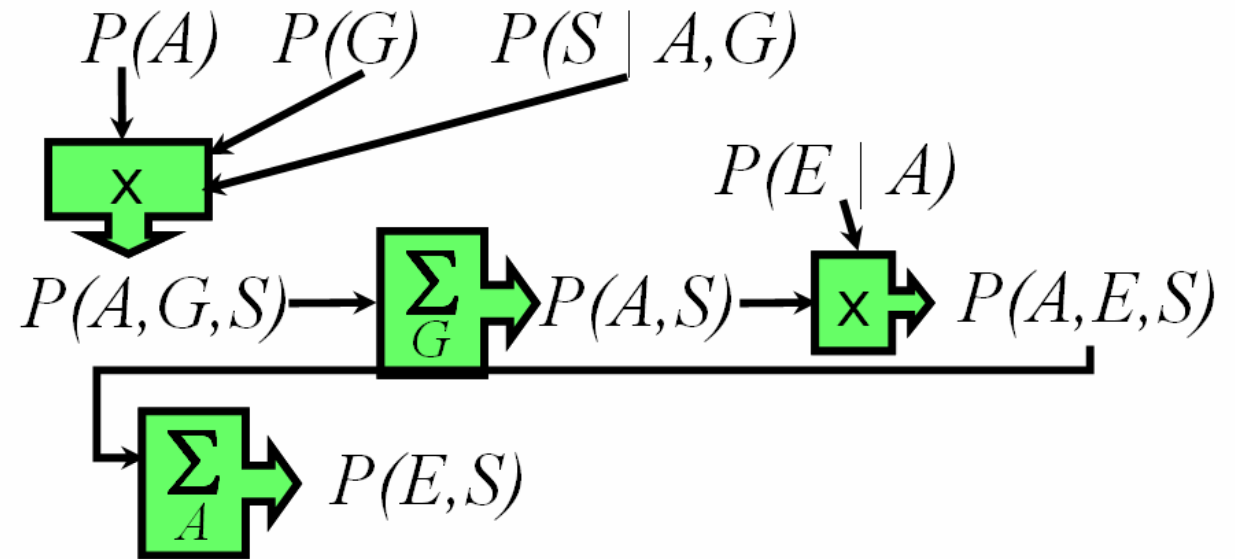
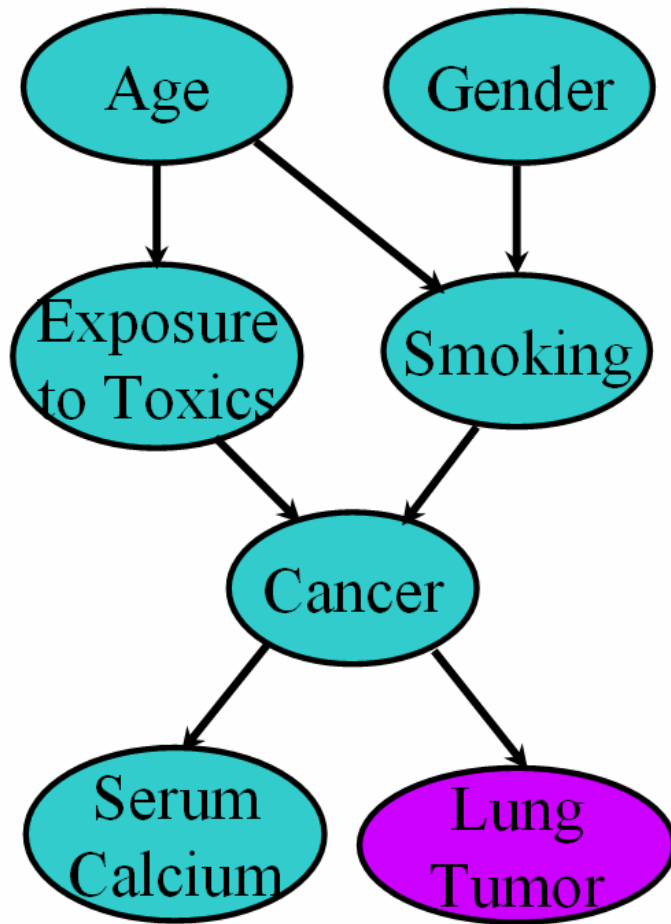
Variable Elimination with loops



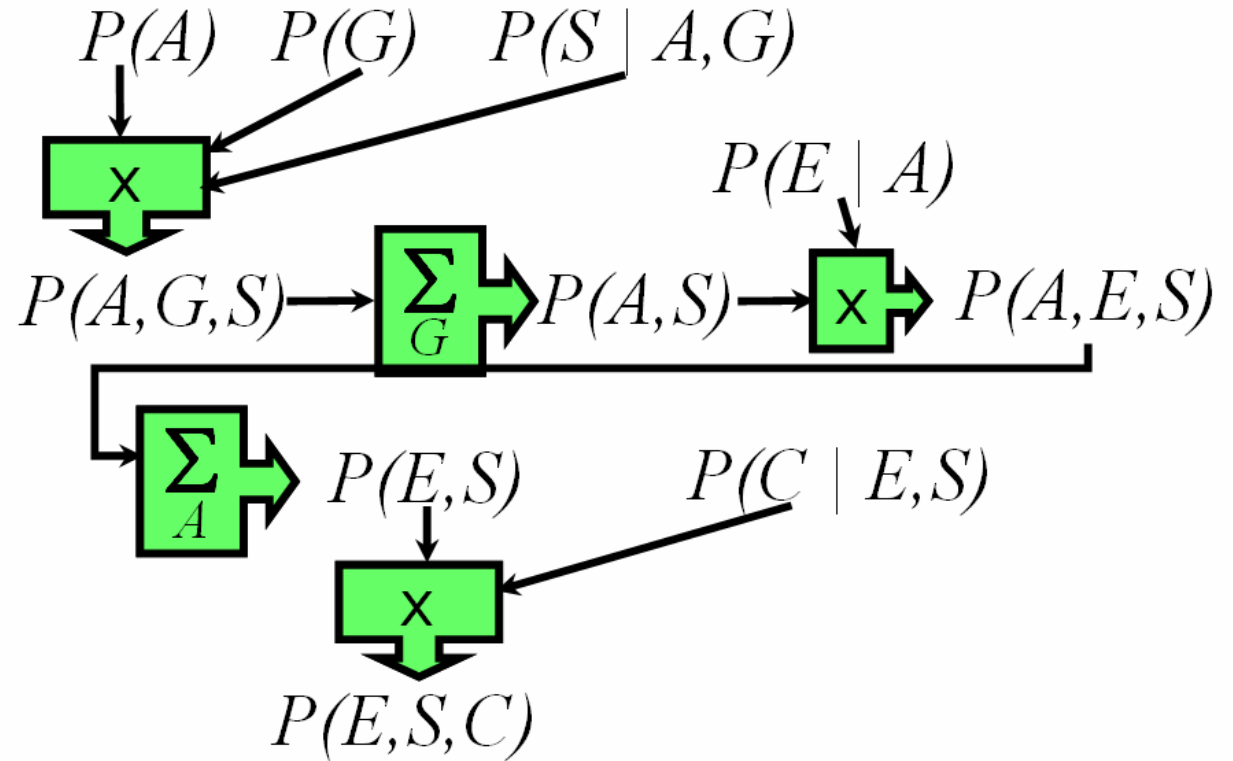
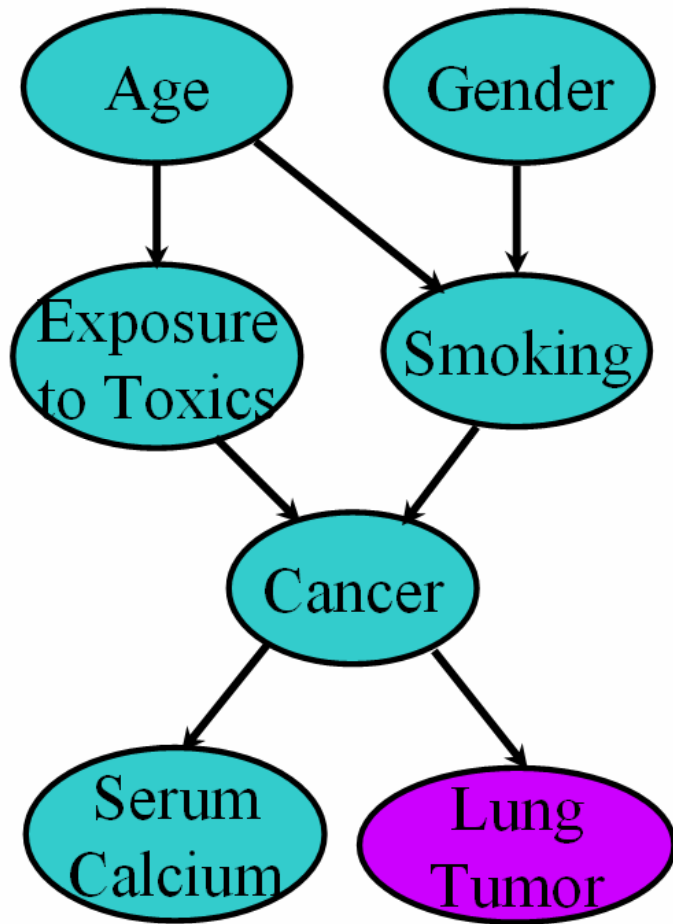
Variable Elimination with loops



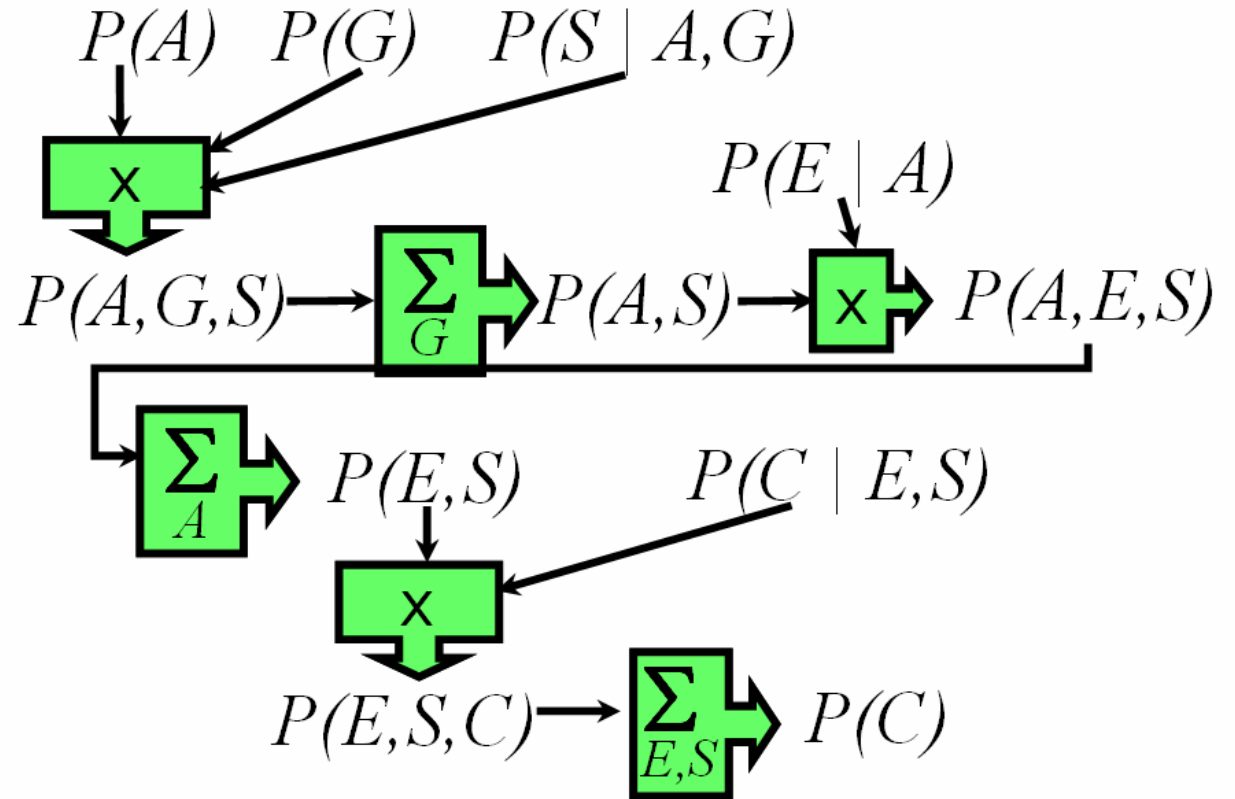
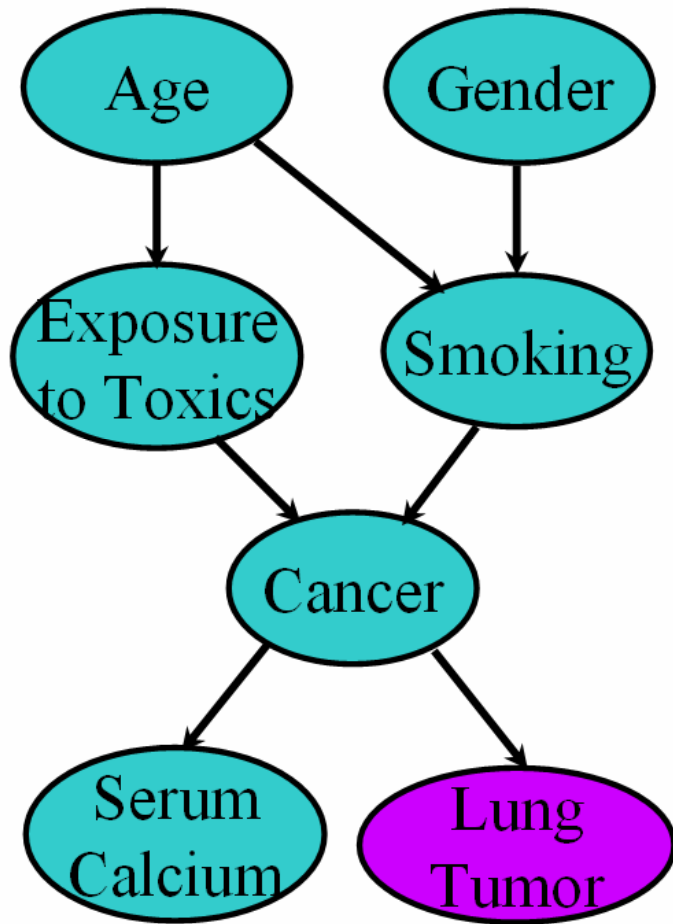
Variable Elimination with loops



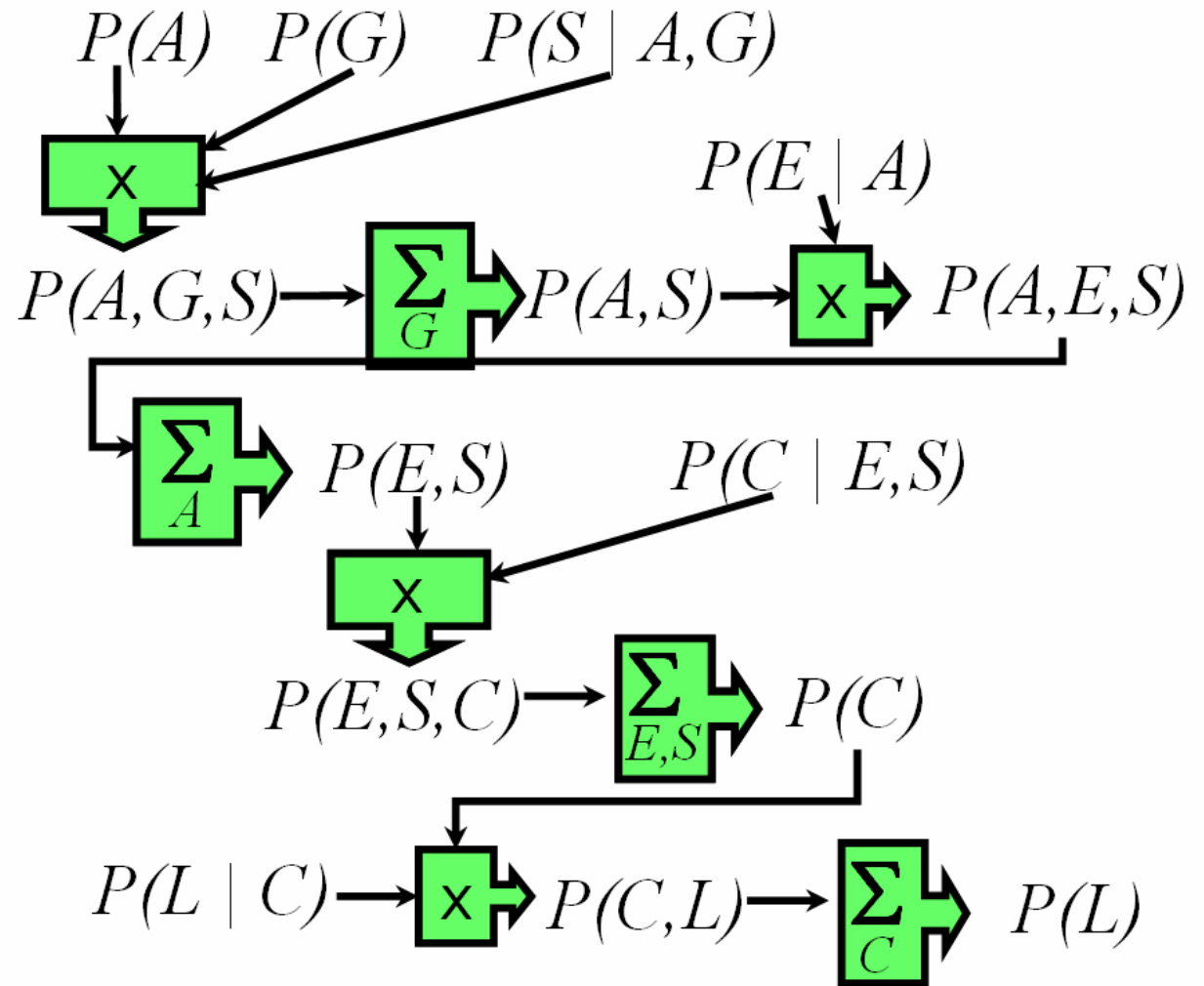
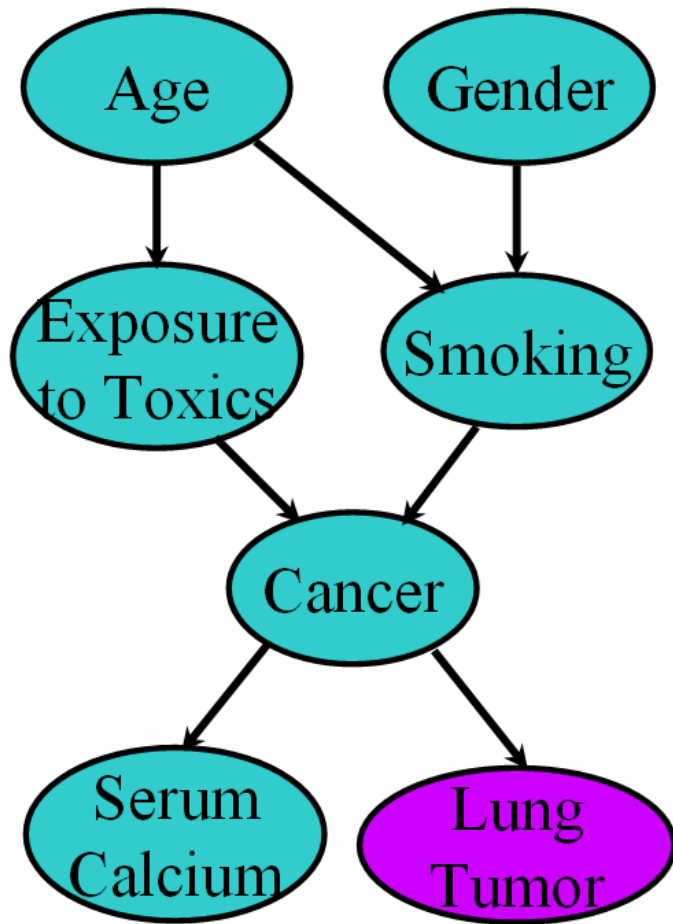
Variable Elimination with loops



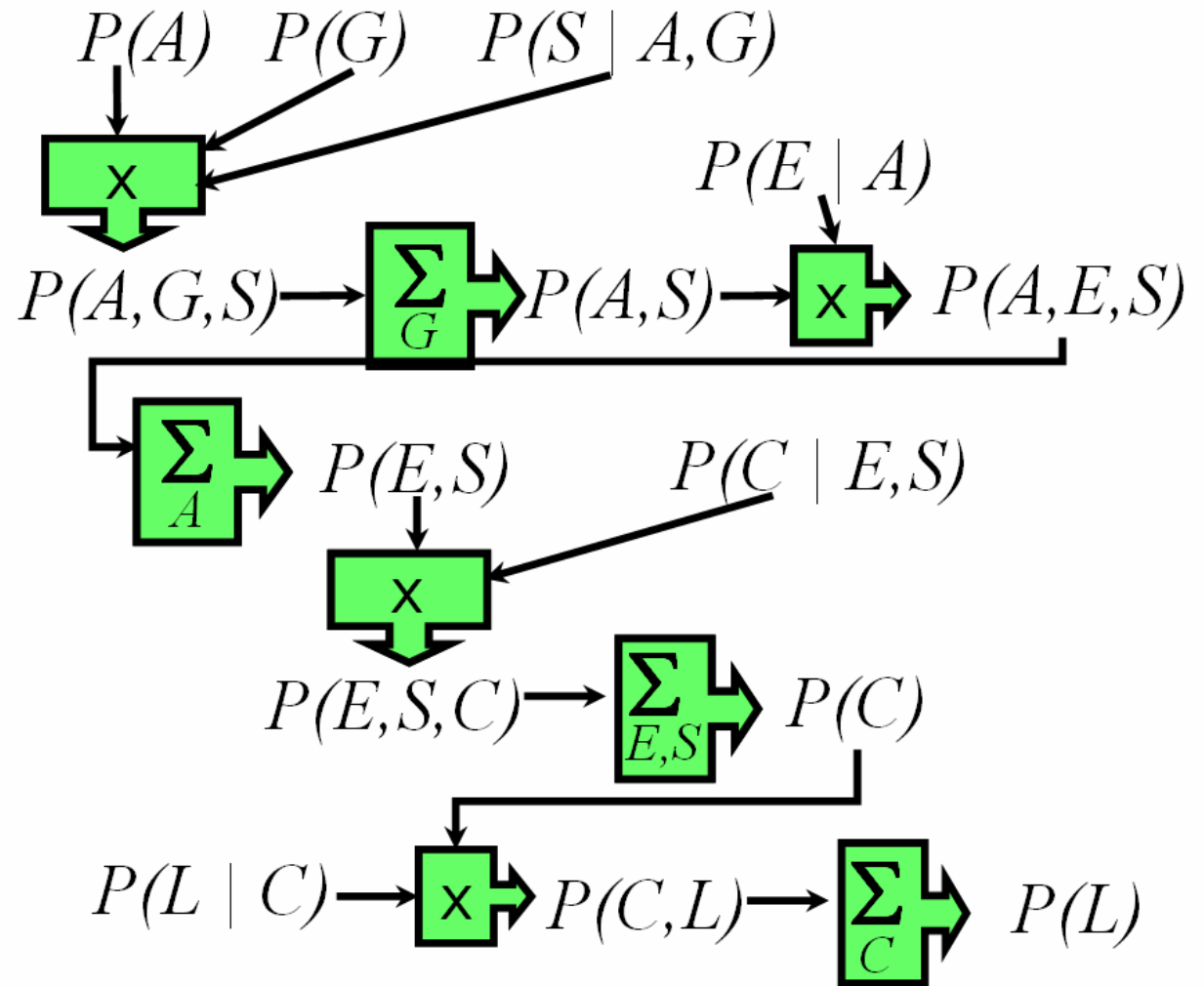
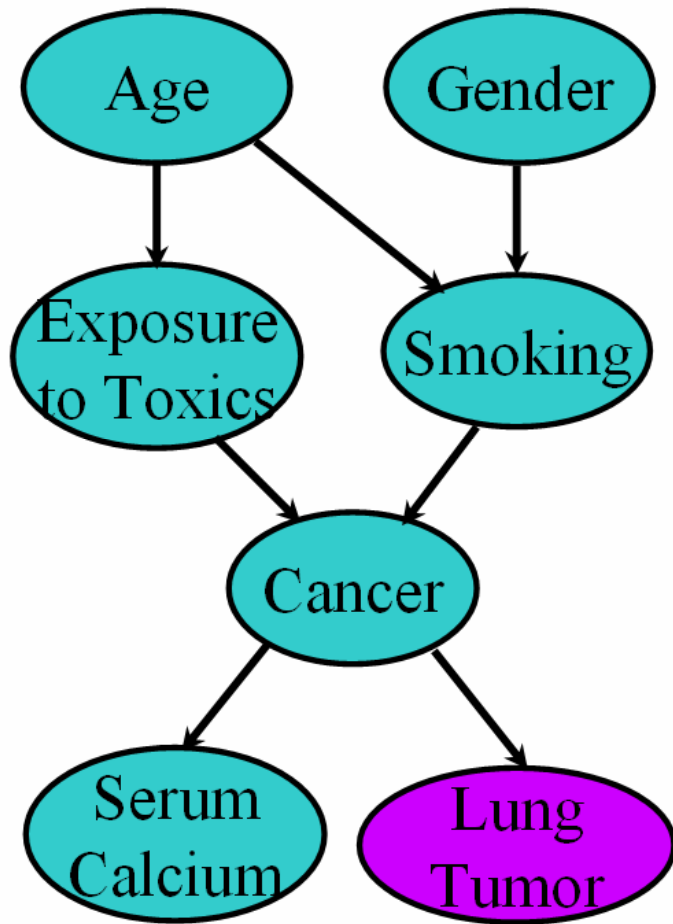
Variable Elimination with loops



Variable Elimination with loops



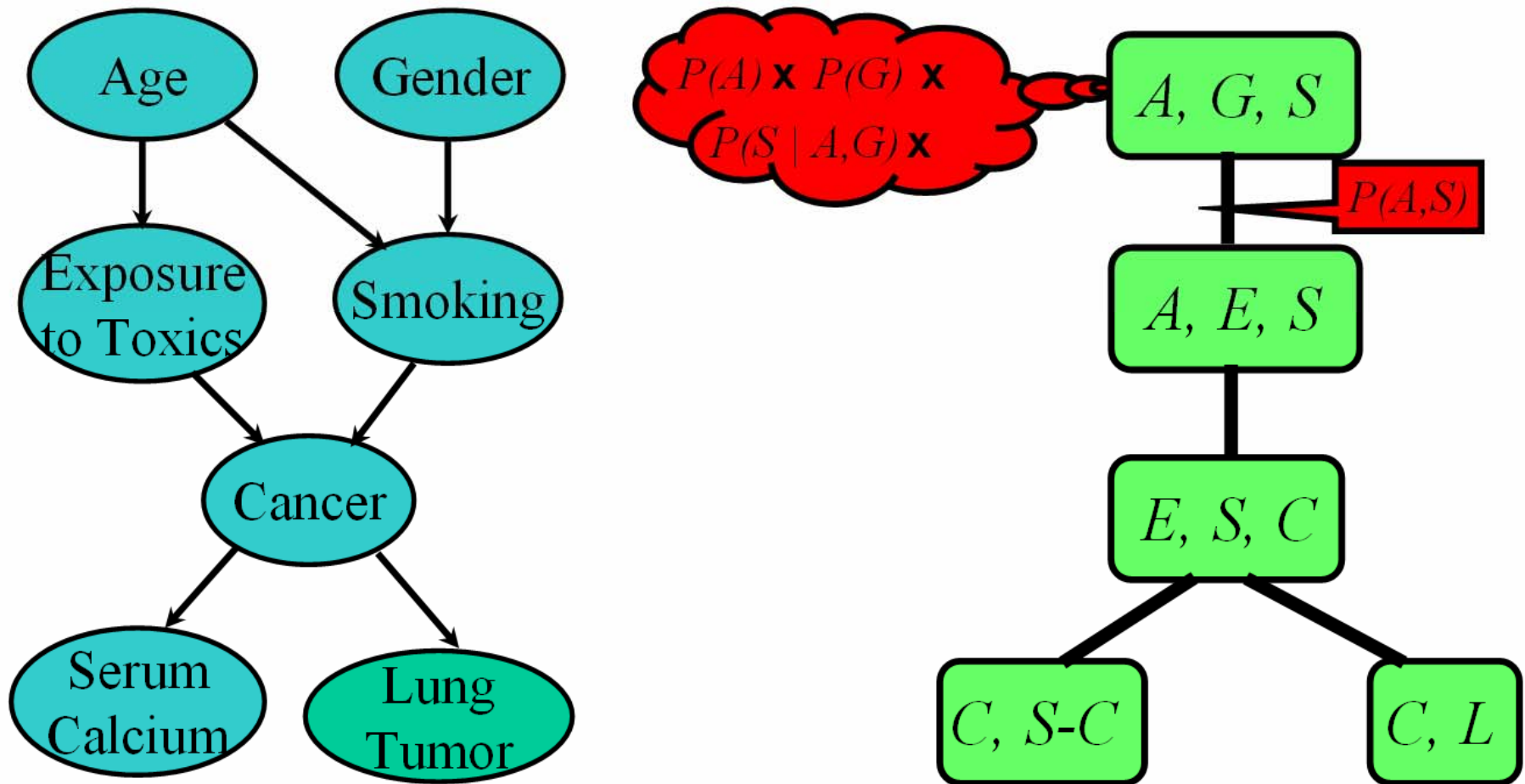
Variable Elimination with loops



Complexity is exponential in the size of the factors

Join trees*

A join tree is a partially precompiled factorization



* aka junction trees, Lauritzen-Spiegelhalter, Hugin alg., ...

Computational complexity

Computational complexity

- **Theorem:** Inference in a multi-connected Bayesian network is NP-hard.

Computational complexity

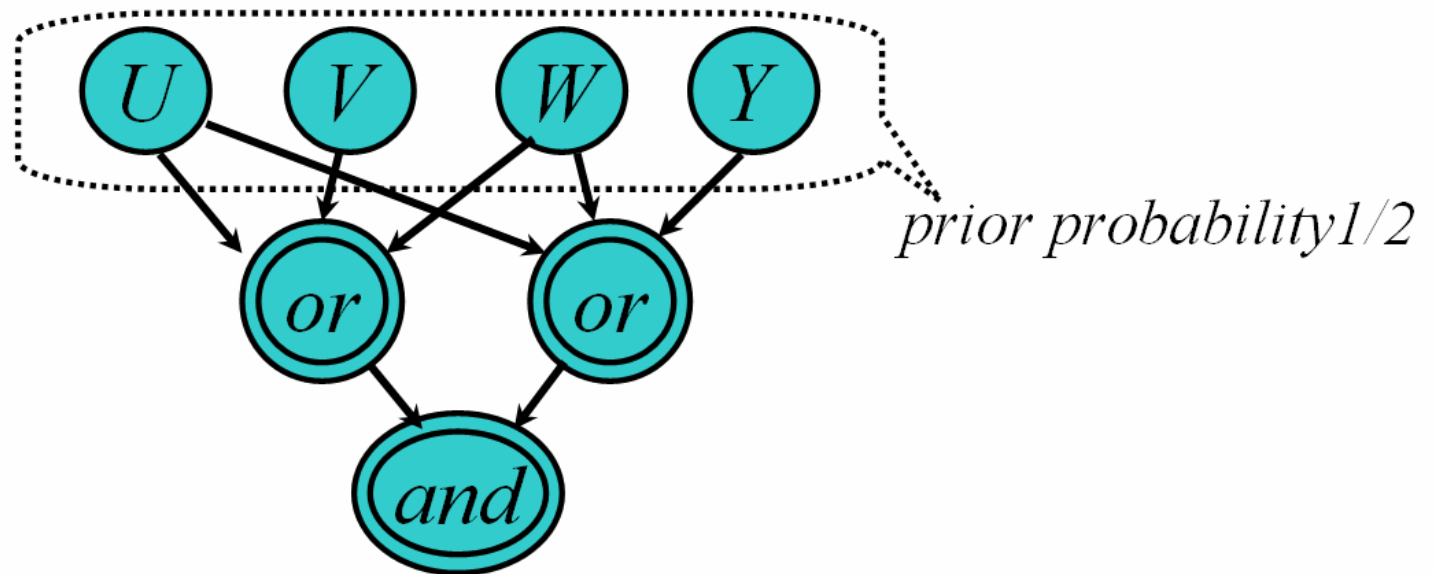
- **Theorem:** Inference in a multi-connected Bayesian network is NP-hard.

Boolean 3CNF formula $\phi = (u \vee \bar{v} \vee w) \wedge (\bar{u} \vee \bar{w} \vee y)$

Computational complexity

- **Theorem:** Inference in a multi-connected Bayesian network is NP-hard.

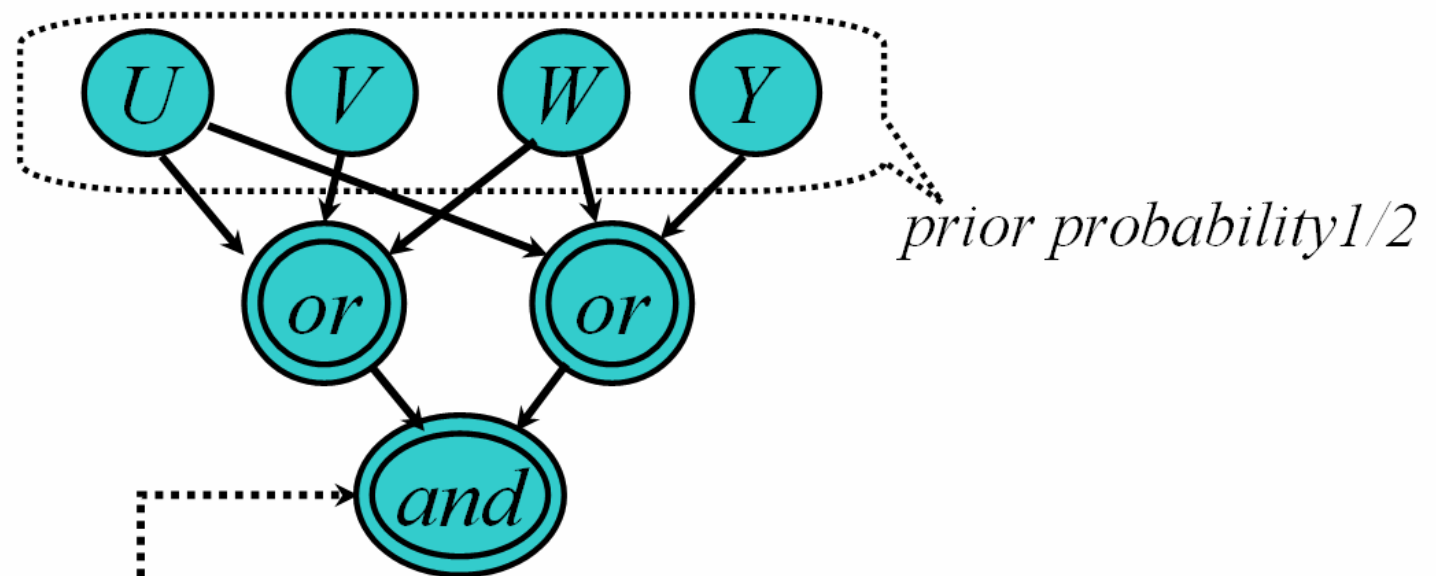
Boolean 3CNF formula $\phi = (u \vee \bar{v} \vee w) \wedge (\bar{u} \vee \bar{w} \vee y)$



Computational complexity

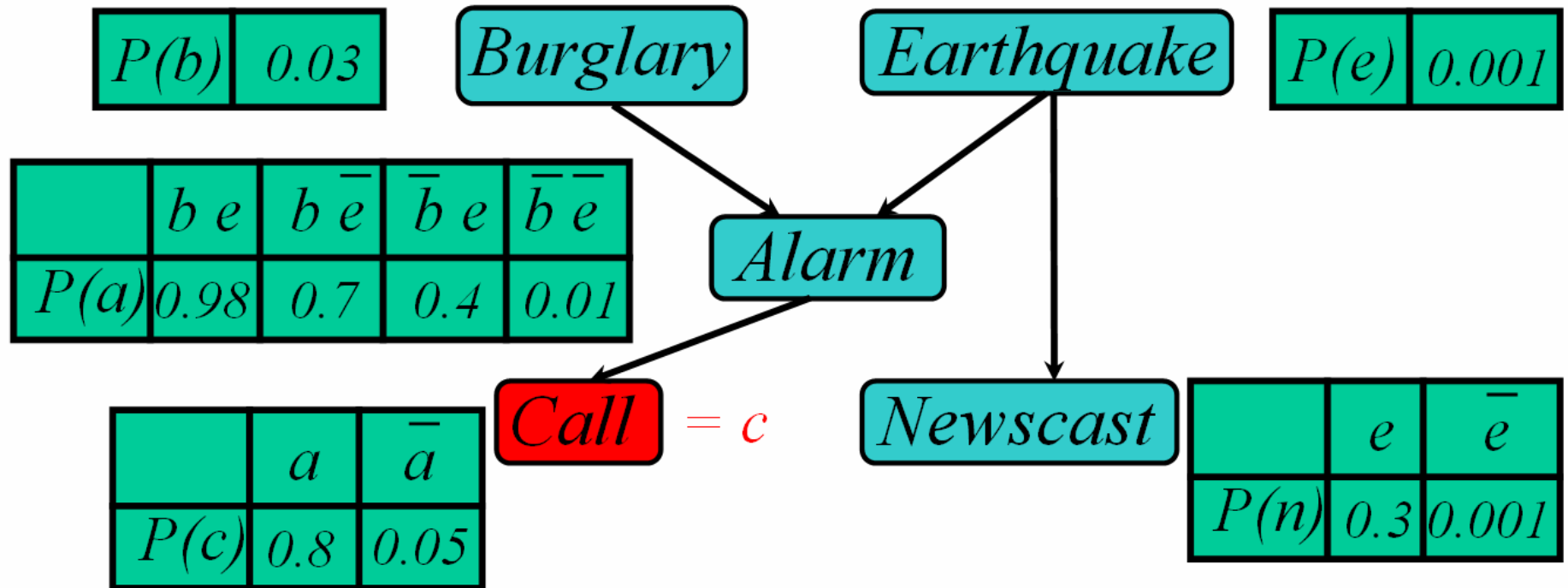
- **Theorem:** Inference in a multi-connected Bayesian network is NP-hard.

Boolean 3CNF formula $\phi = (u \vee \bar{v} \vee w) \wedge (\bar{u} \vee \bar{w} \vee y)$

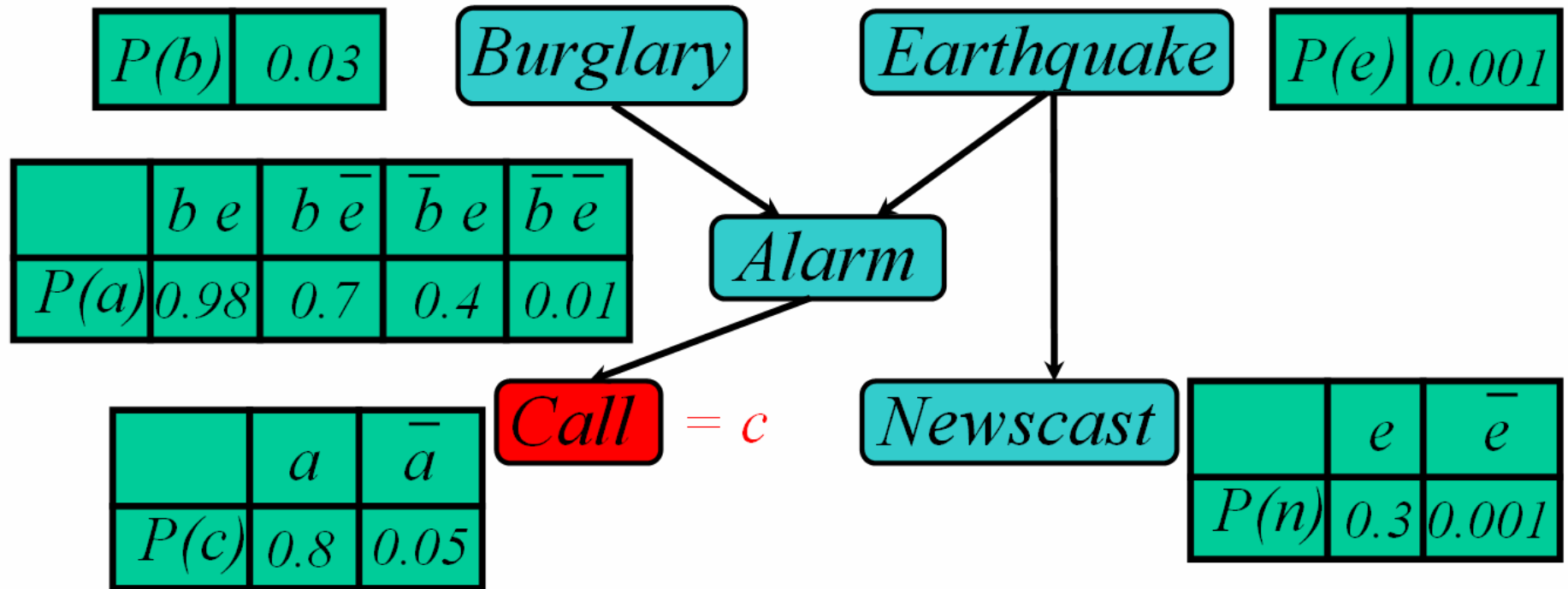


$Probability(\cdot) = 1/2^n \cdot \# \text{ satisfying assignments of } \phi$

Stochastic simulation



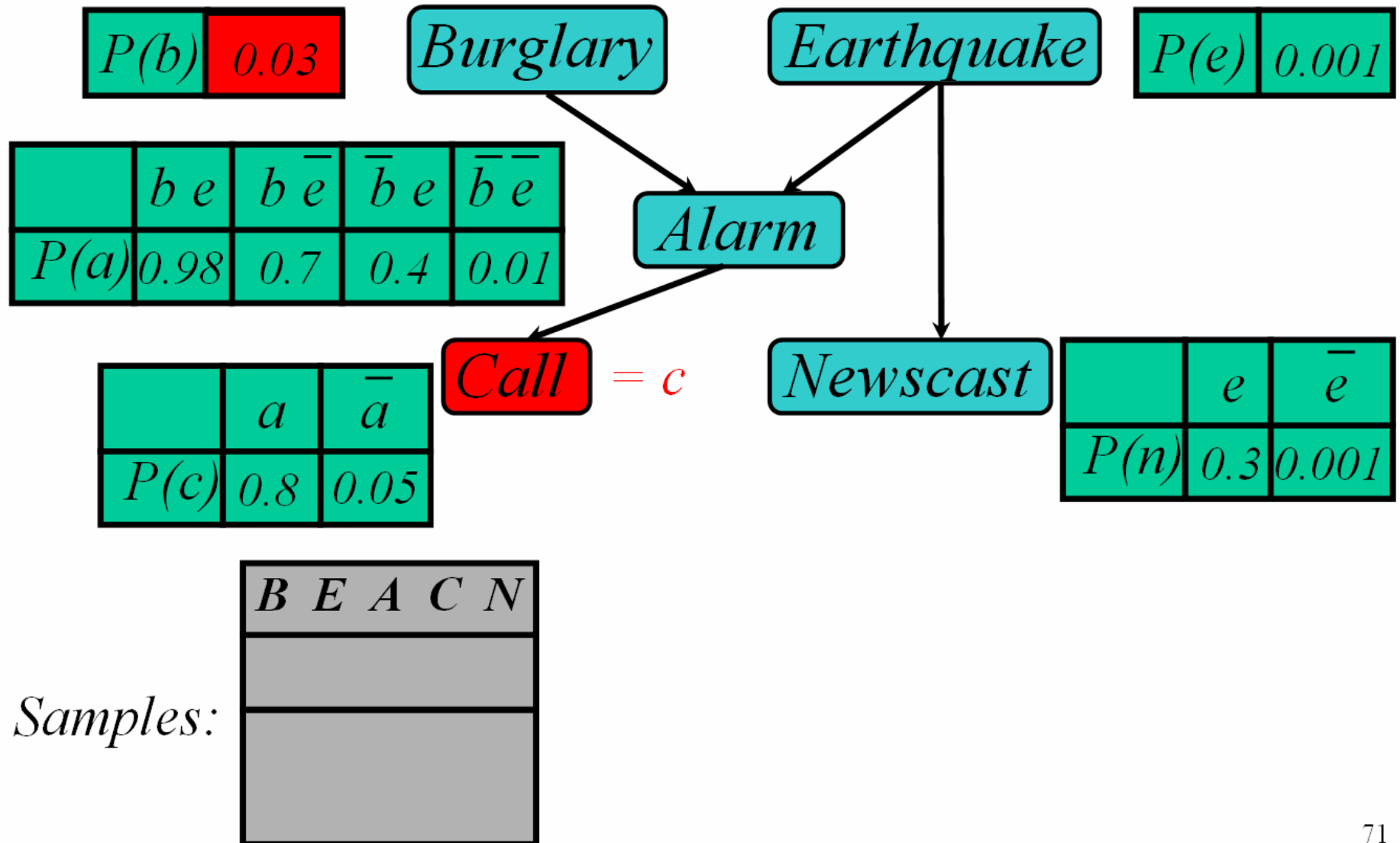
Stochastic simulation



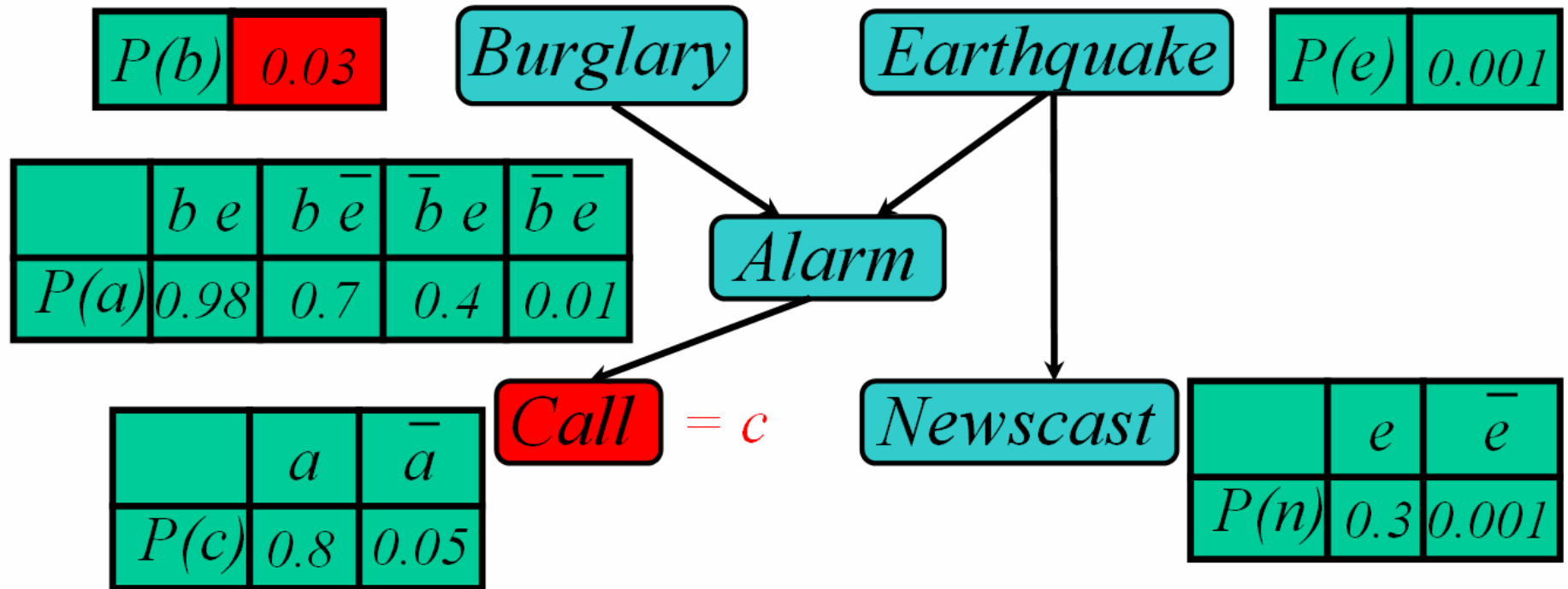
Samples:

<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>

Stochastic simulation



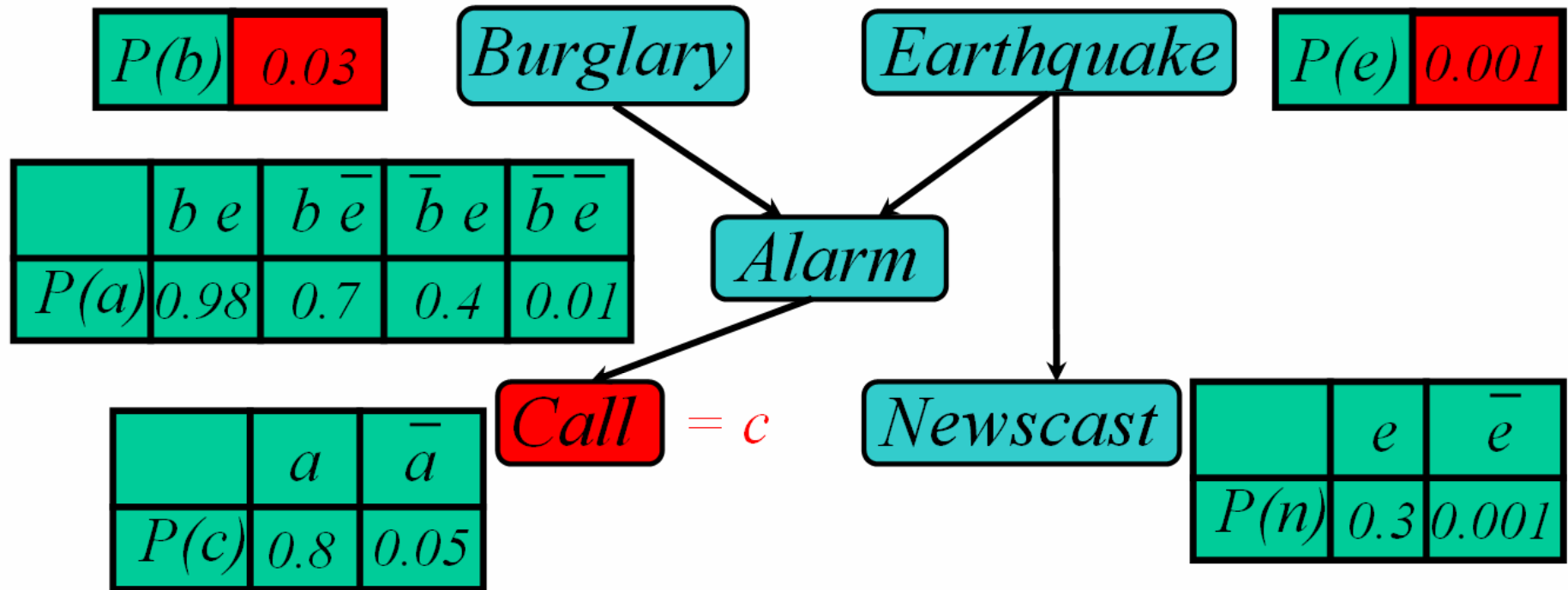
Stochastic simulation



Samples:

B	E	A	C	N
\bar{b}				

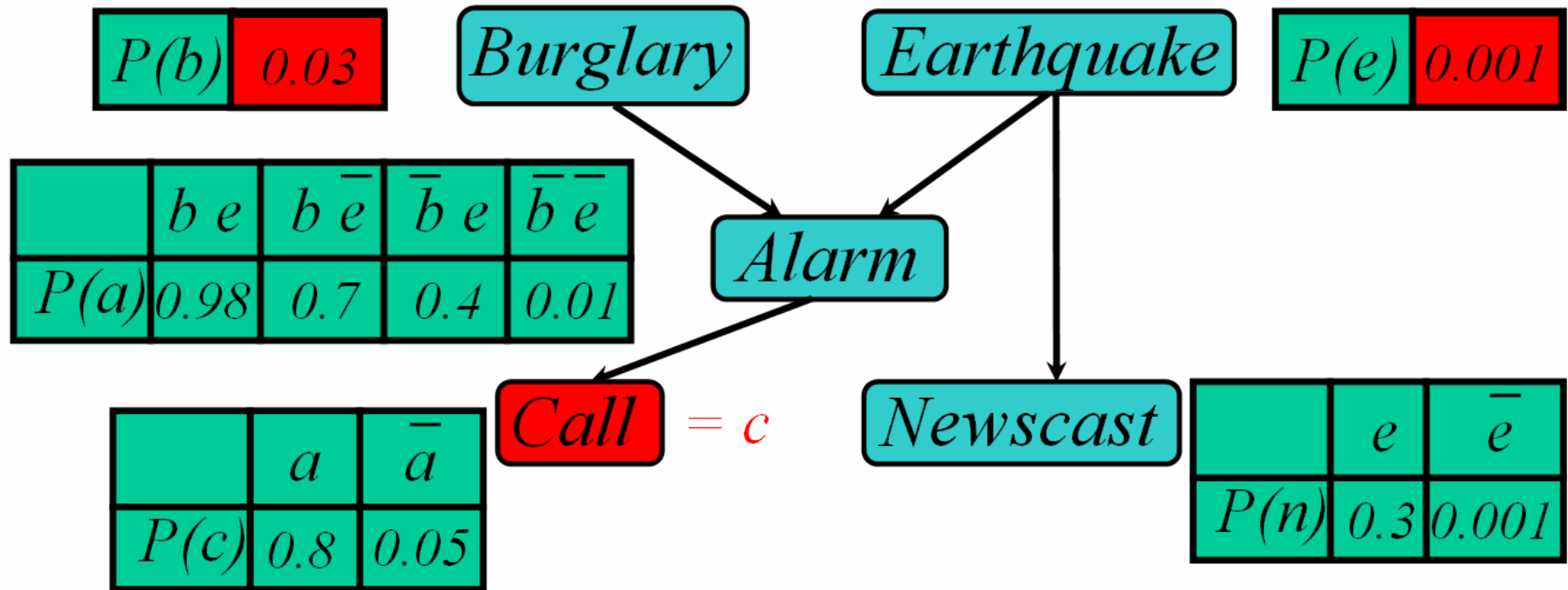
Stochastic simulation



Samples:

<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}				

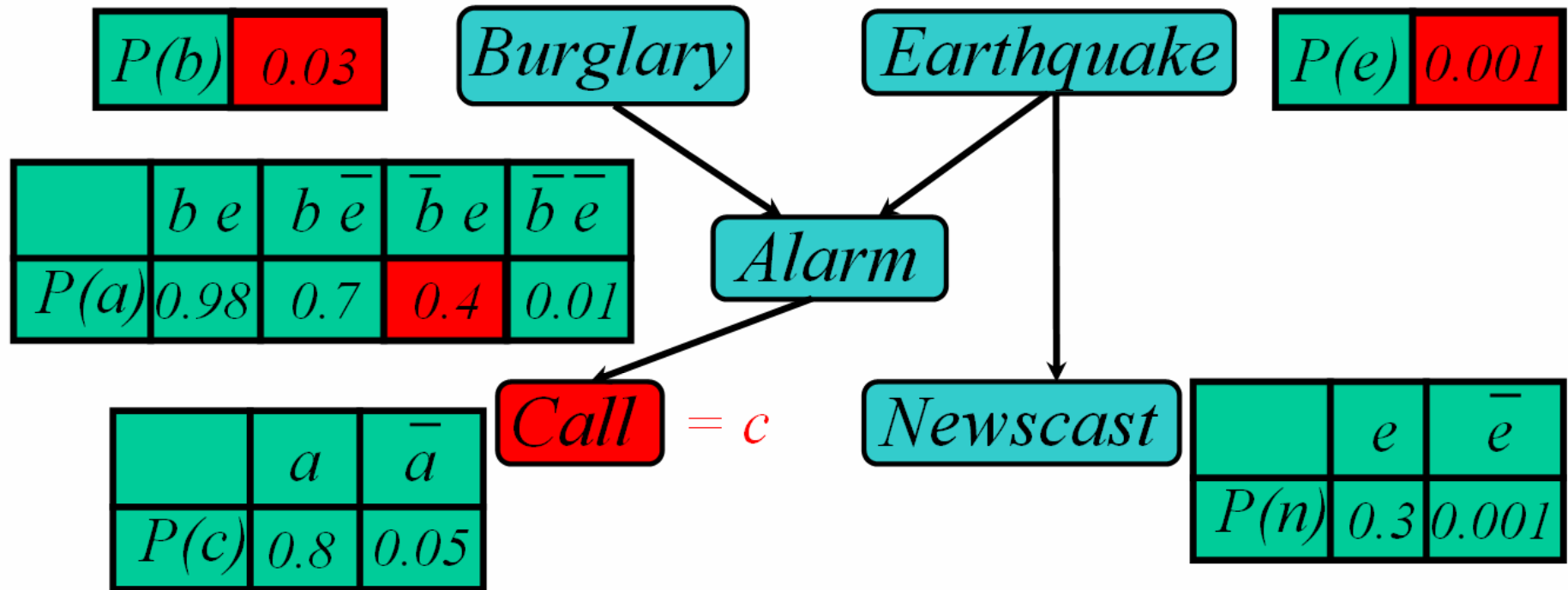
Stochastic simulation



Samples:

<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}	e			

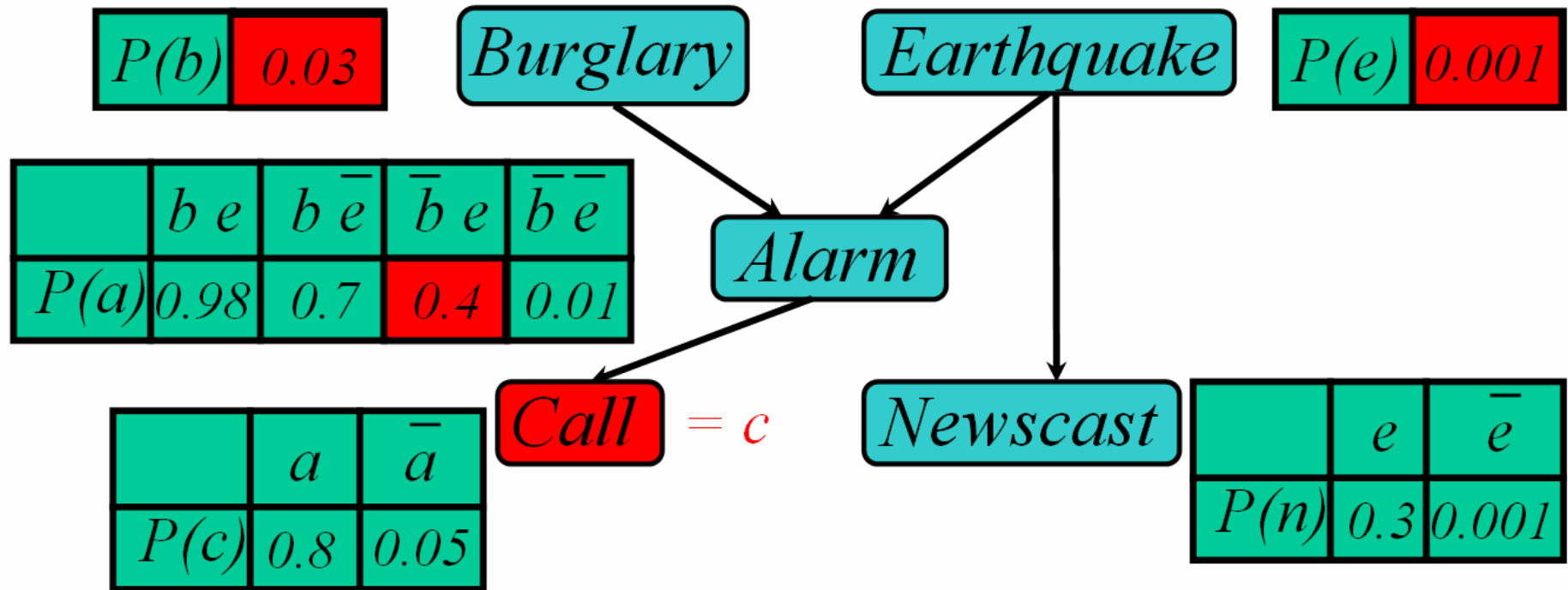
Stochastic simulation



Samples:

B	E	A	C	N
\bar{b}	e			

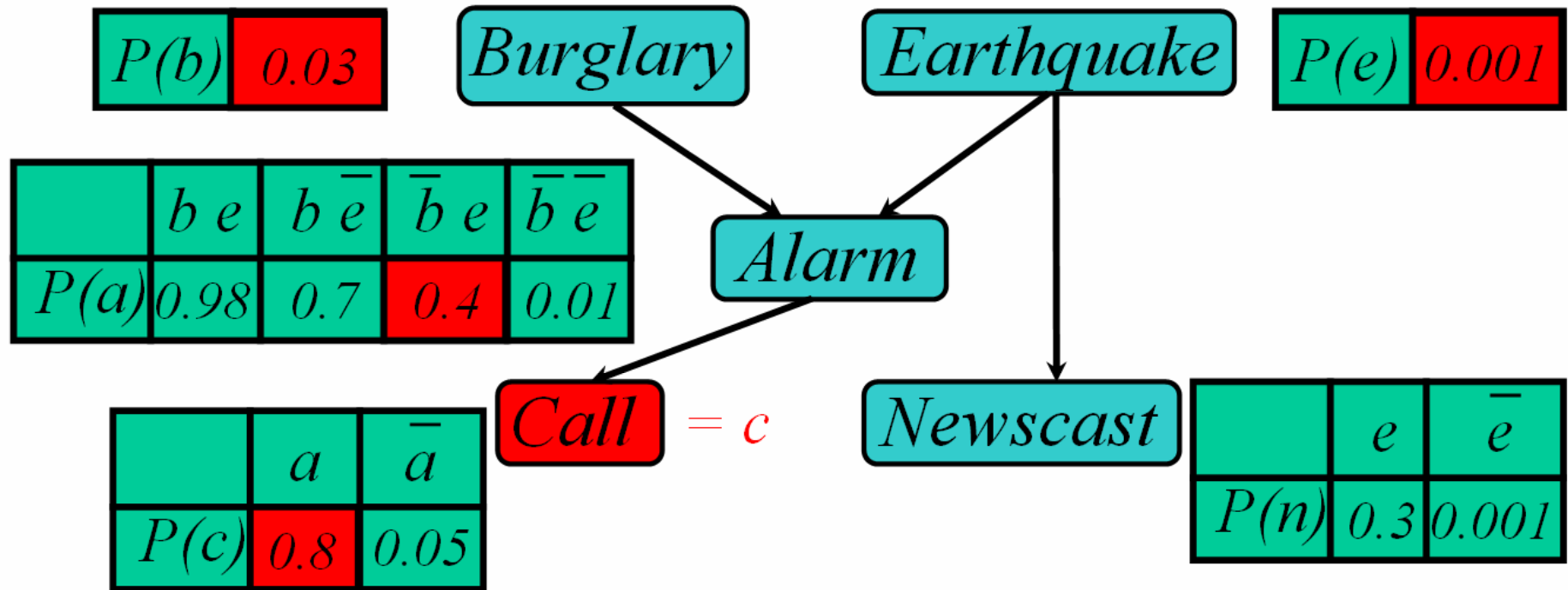
Stochastic simulation



Samples:

B	E	A	C	N
\bar{b}	e	a		

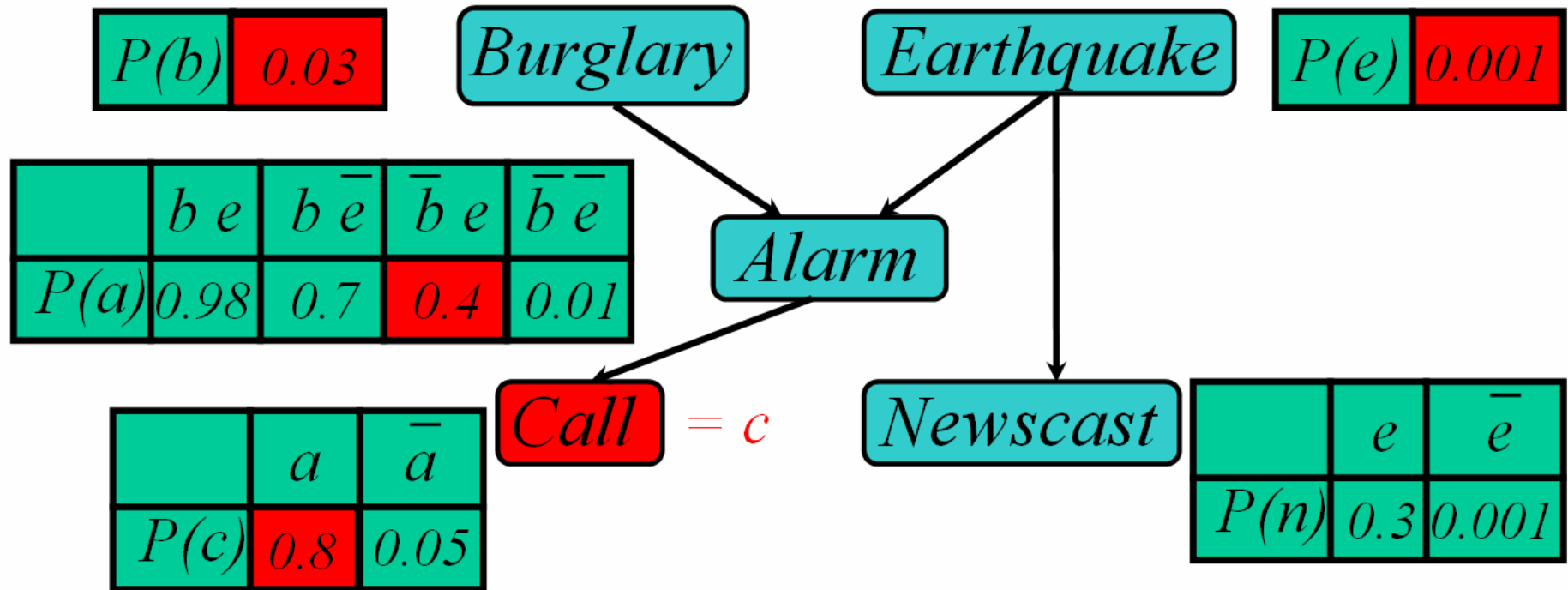
Stochastic simulation



Samples:

B	E	A	C	N
\bar{b}	e	a		

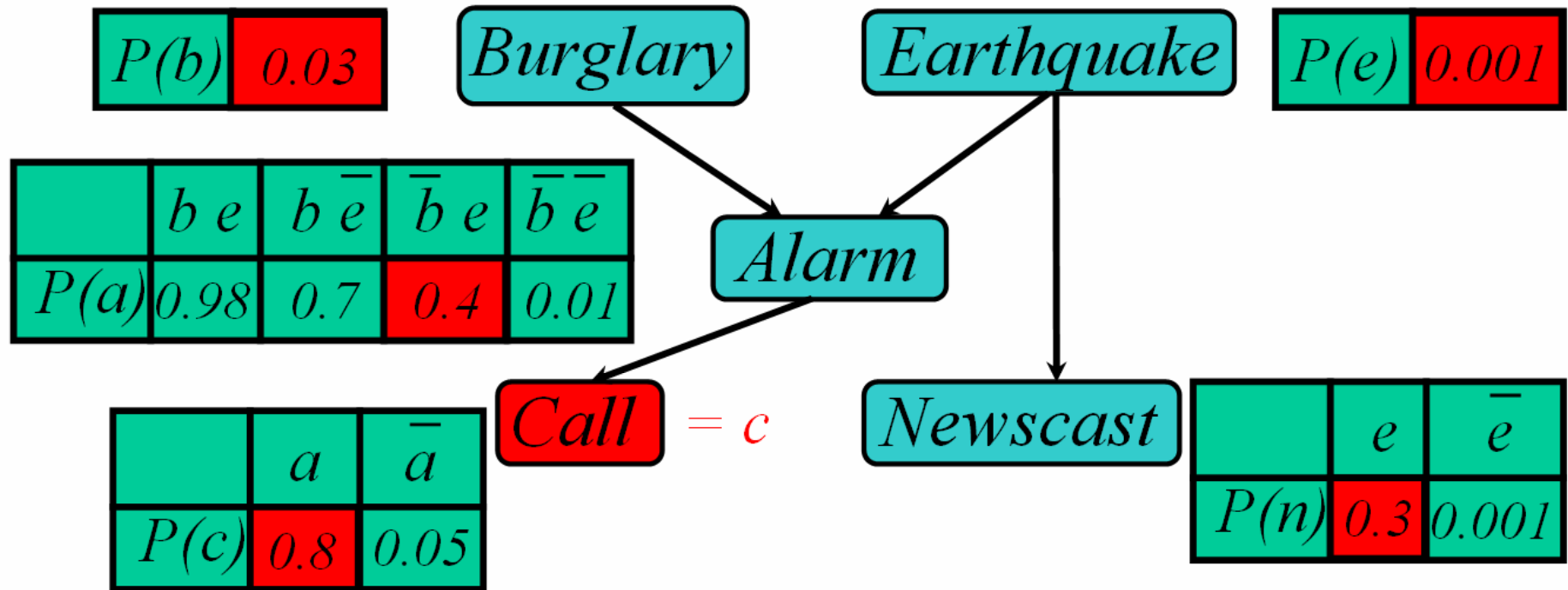
Stochastic simulation



Samples:

B	E	A	C	N
\bar{b}	e	a	c	

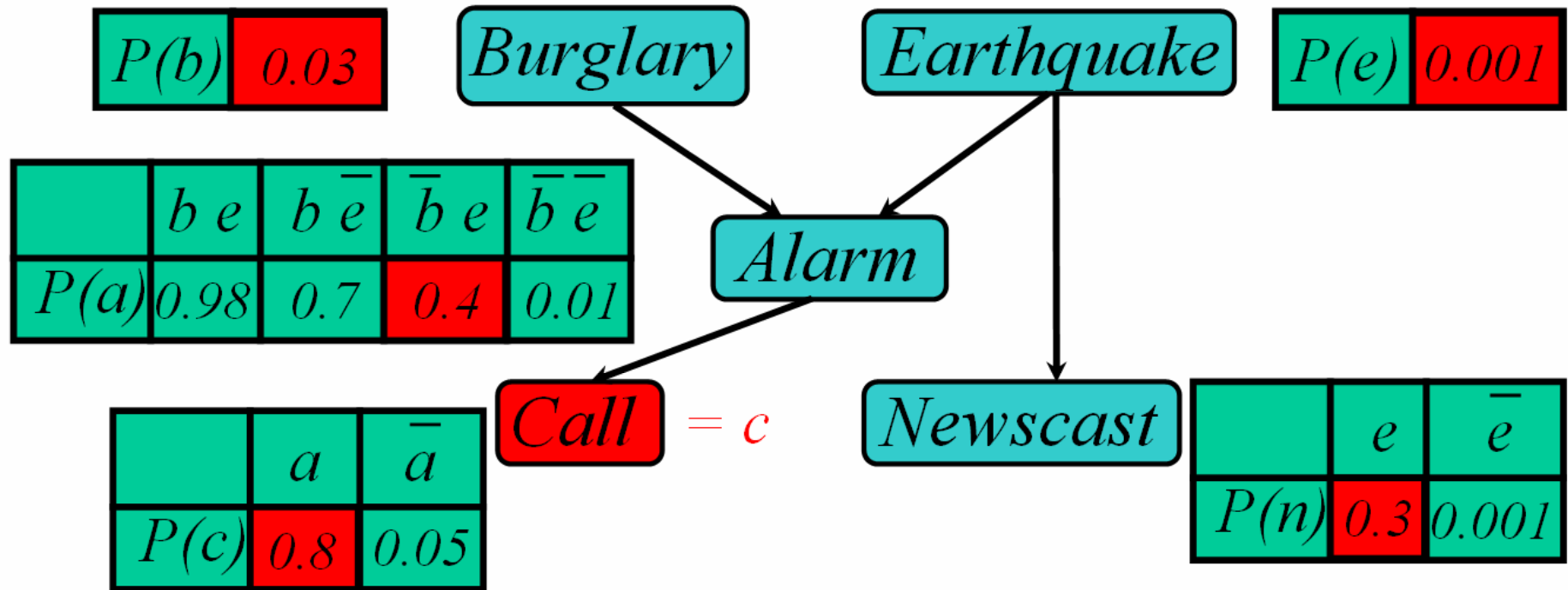
Stochastic simulation



Samples:

B	E	A	C	N
\bar{b}	e	a	c	

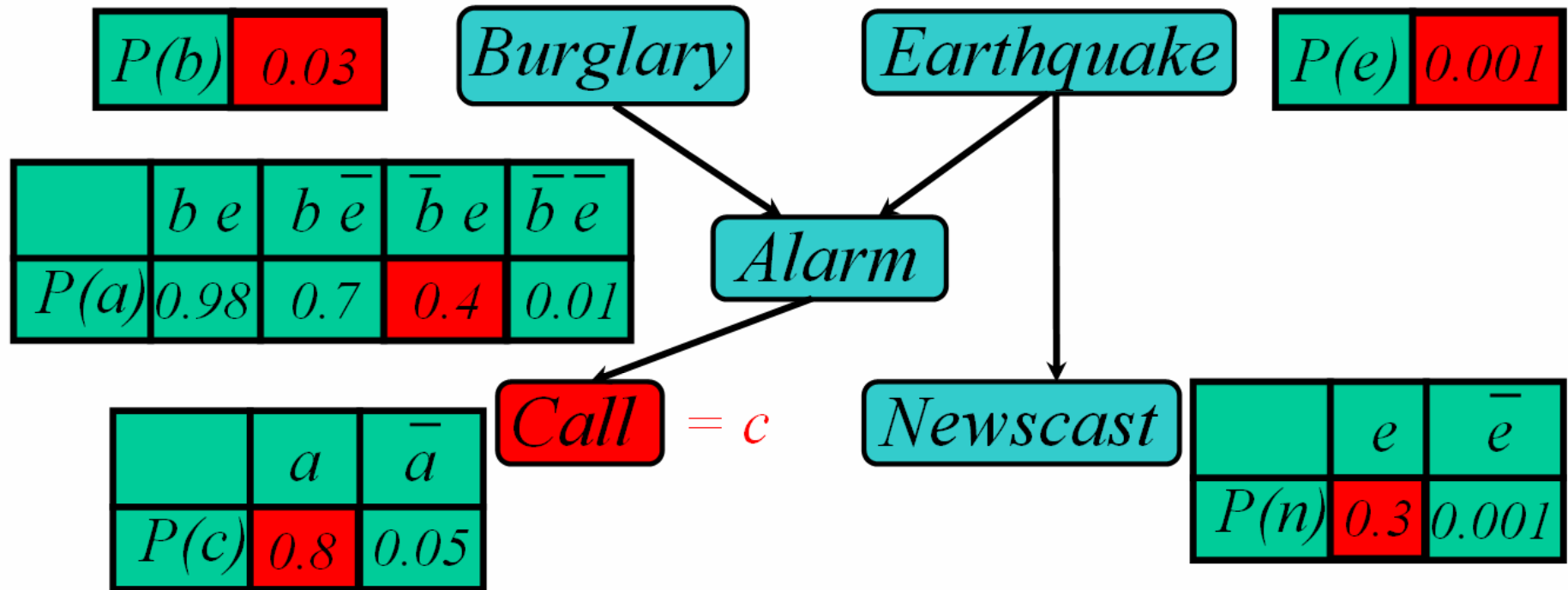
Stochastic simulation



Samples:

B	E	A	C	N
\bar{b}	e	a	c	\bar{n}

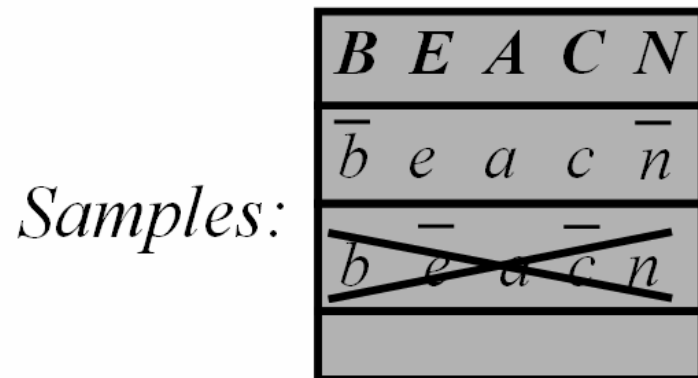
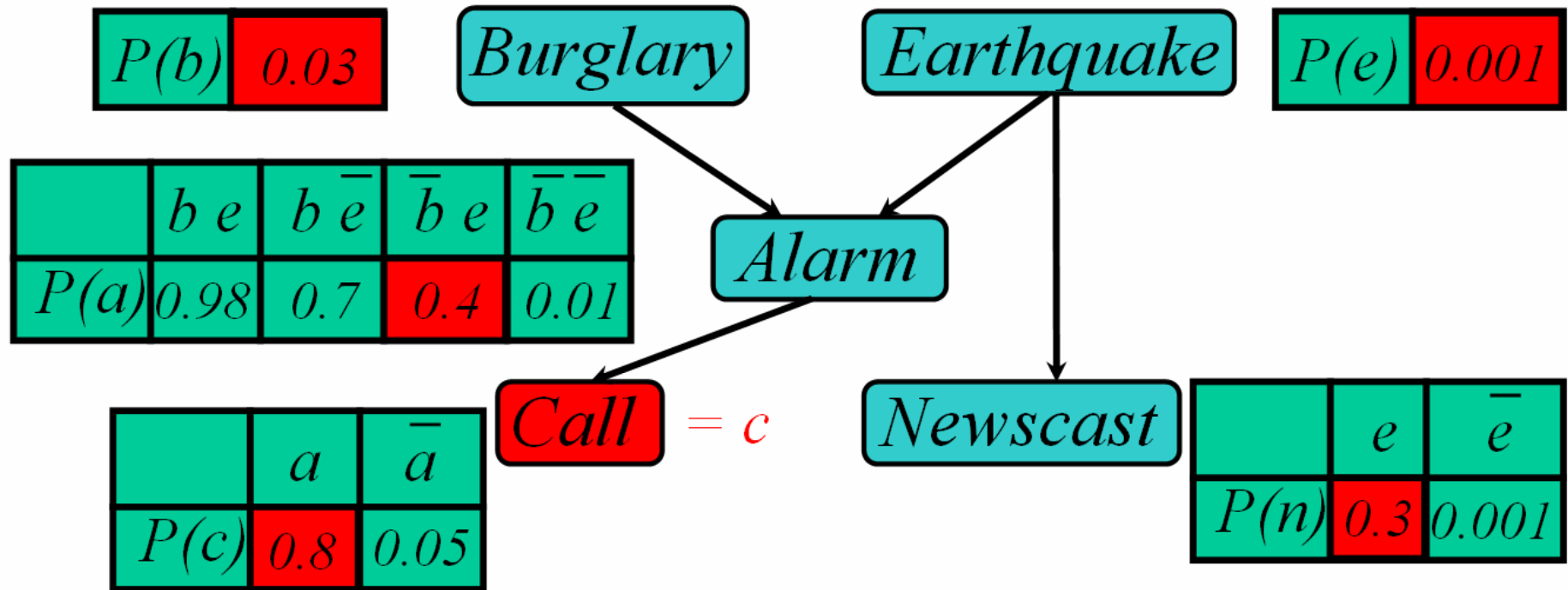
Stochastic simulation



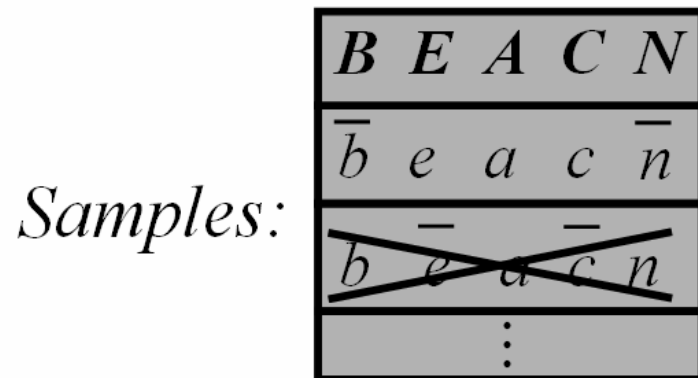
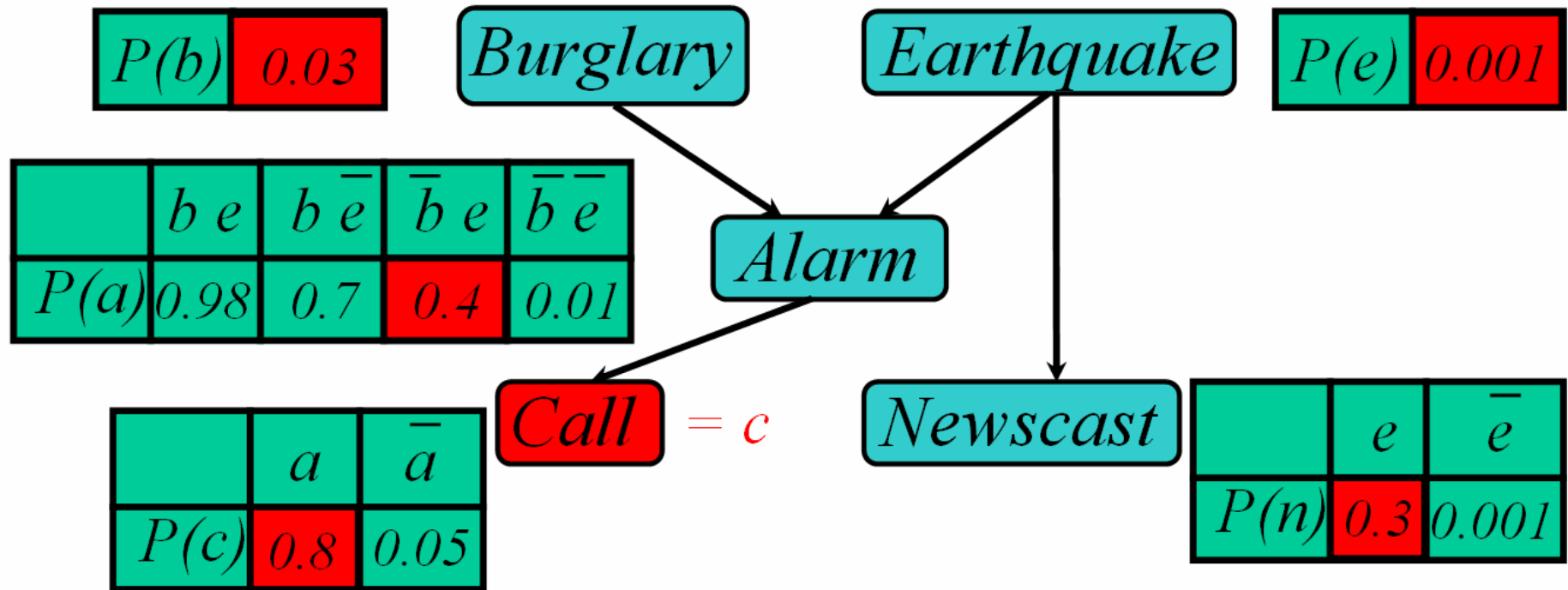
Samples:

<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}	e	a	c	\bar{n}
b	\bar{e}	a	\bar{c}	n

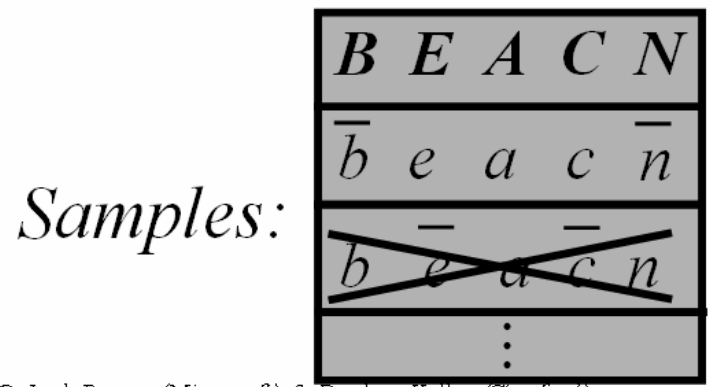
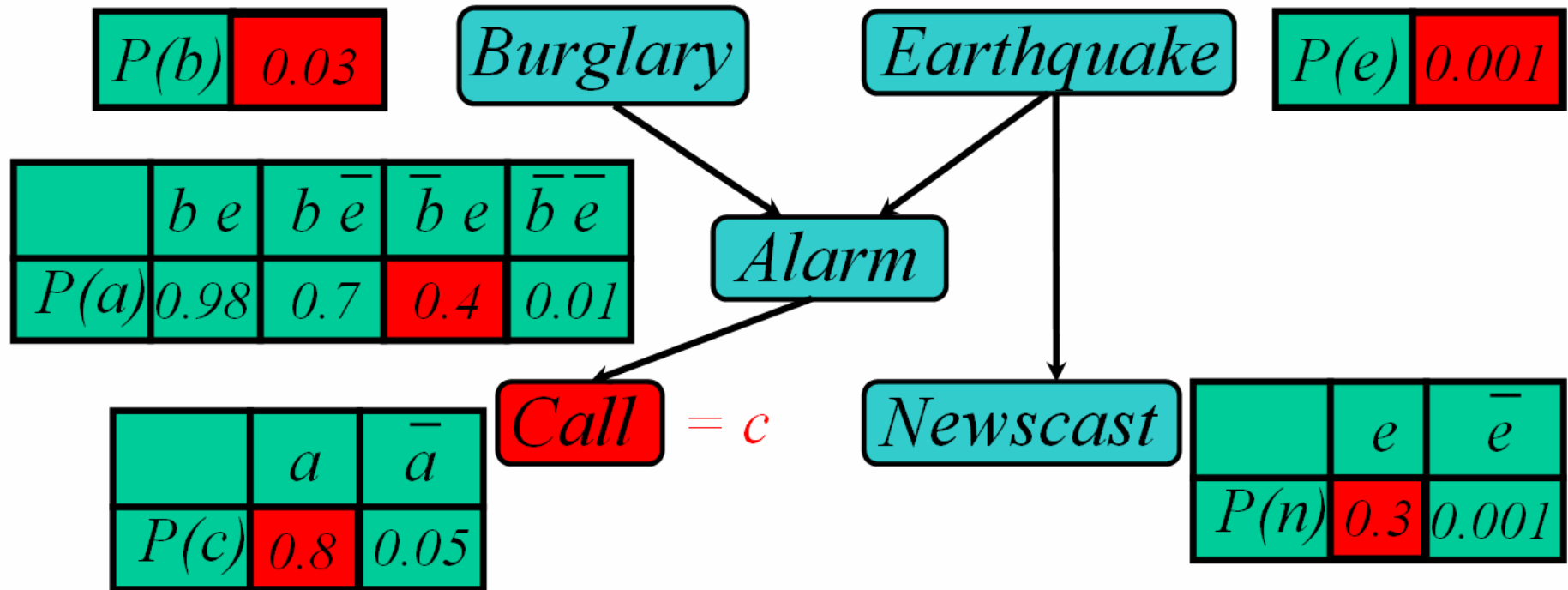
Stochastic simulation



Stochastic simulation

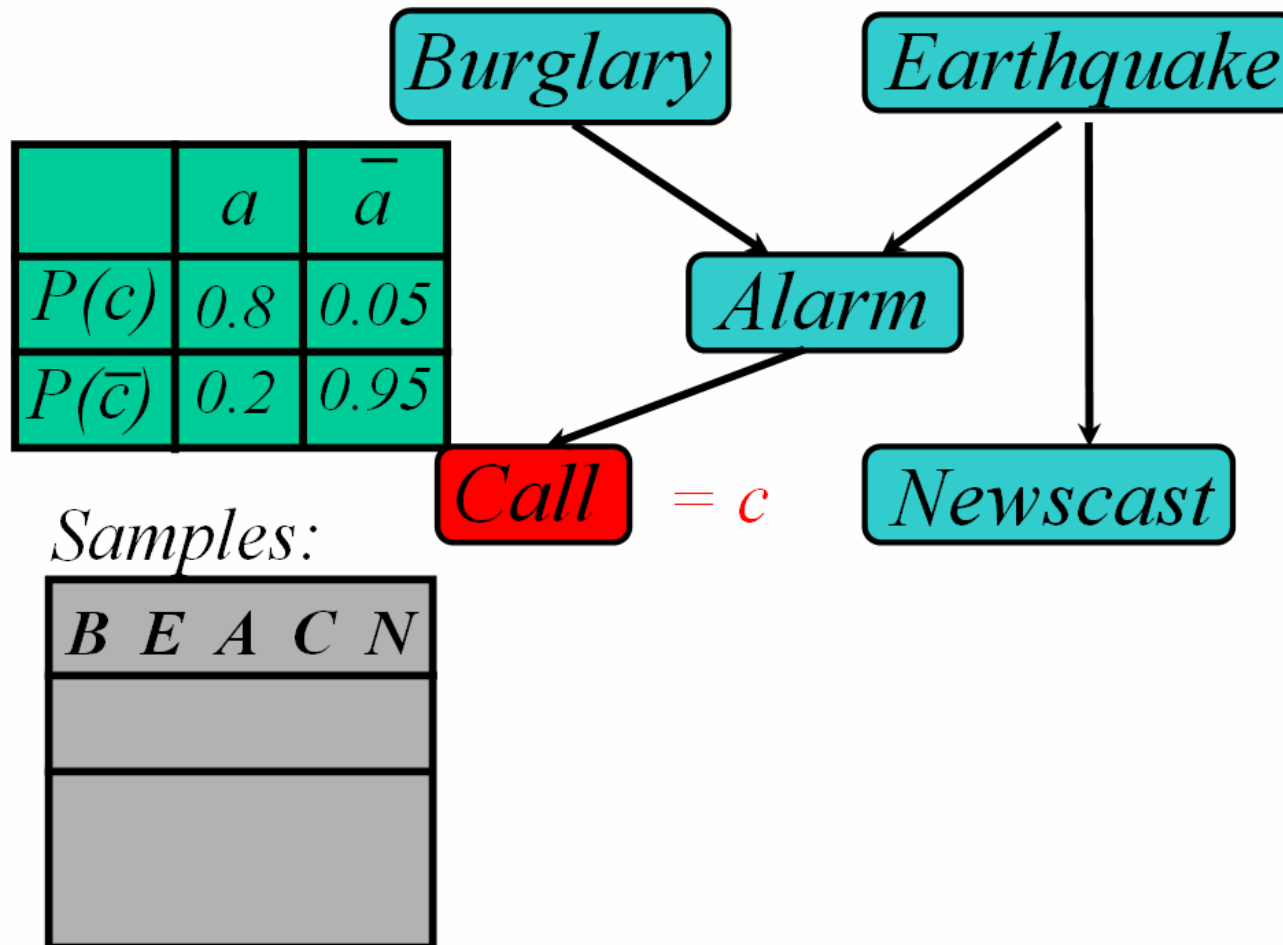


Stochastic simulation

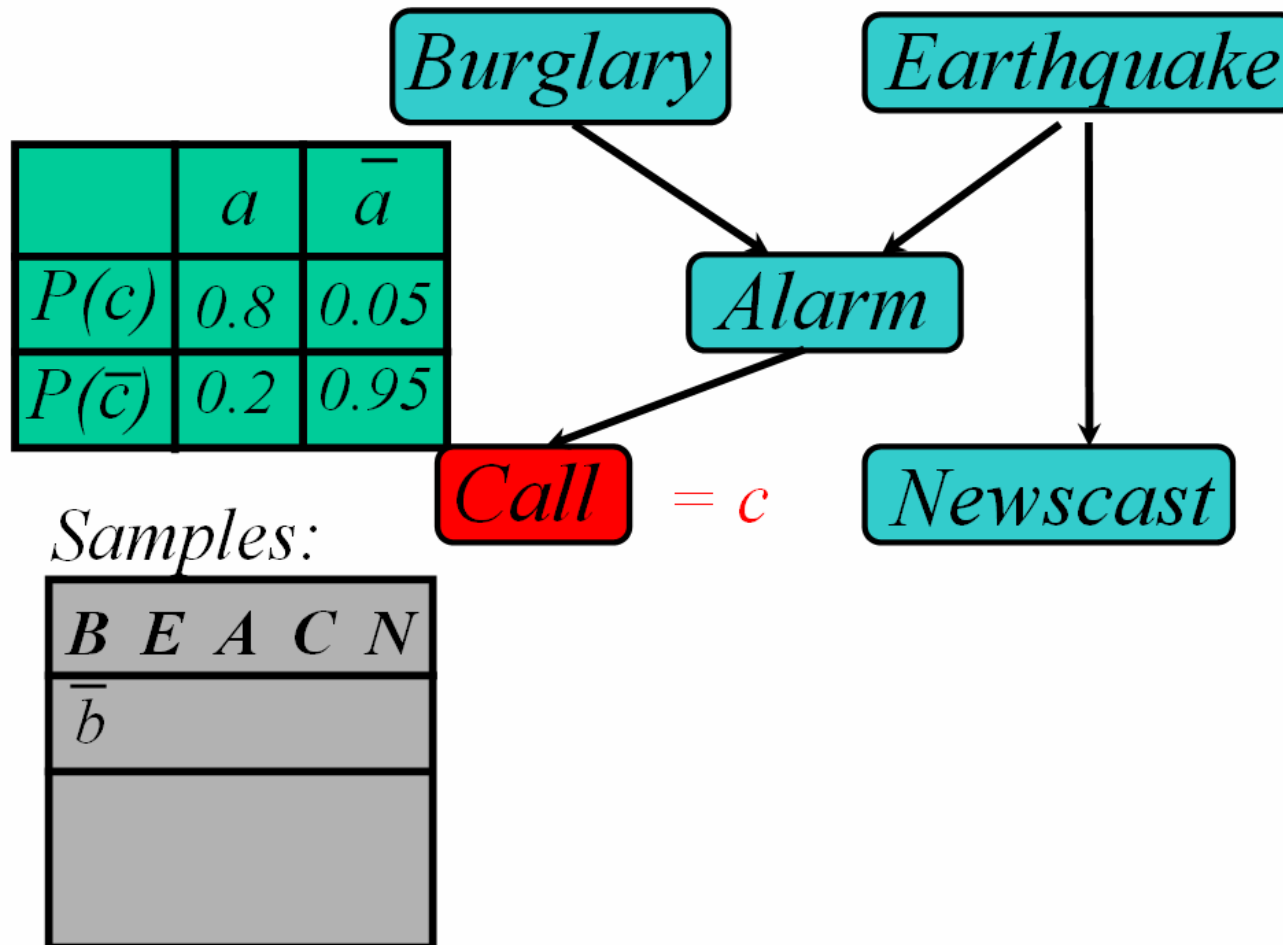


$$P(b|c) \sim \frac{\text{\# of live samples with } B=b}{\text{total \# of live samples}}$$

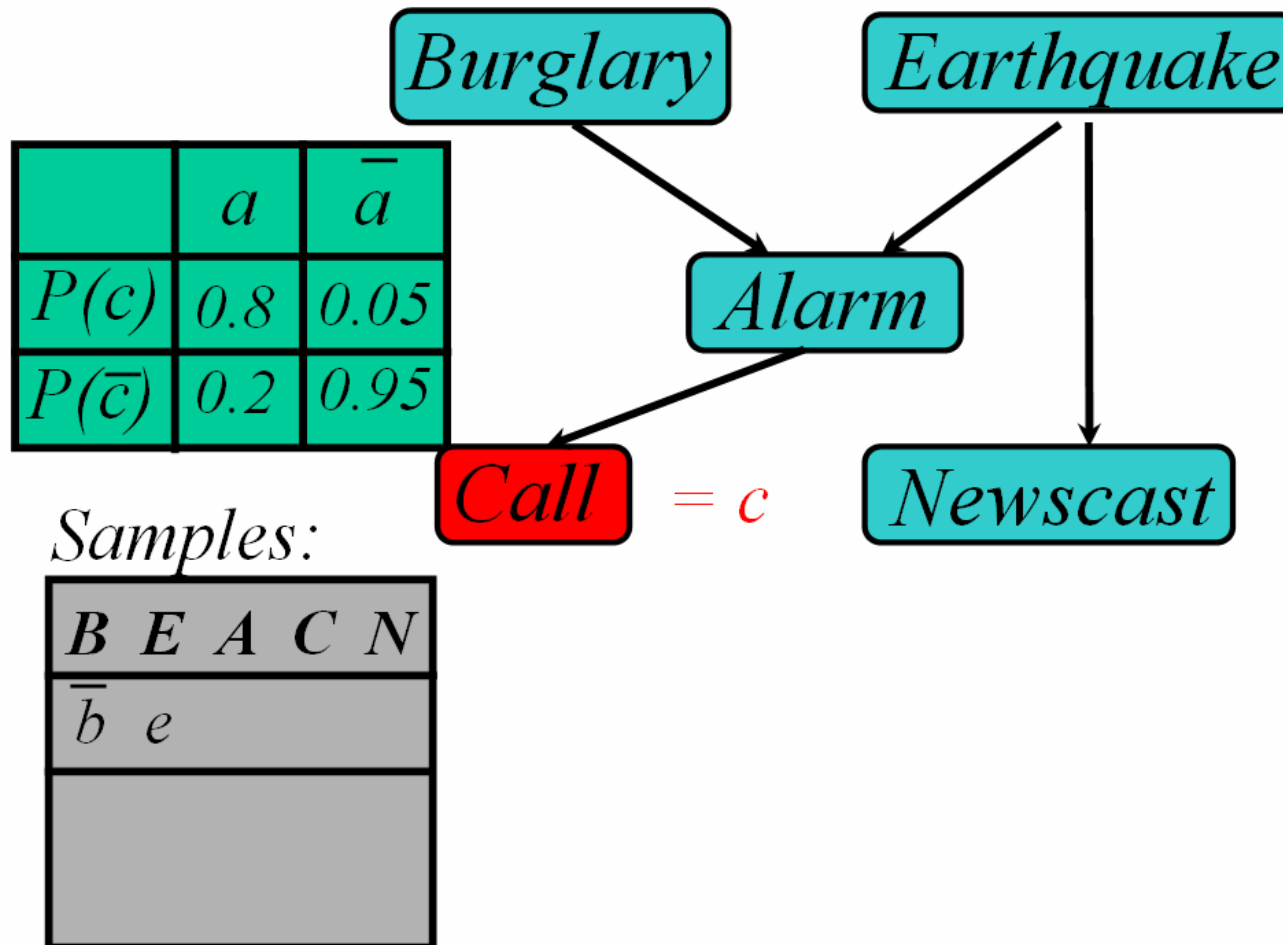
Likelihood weighting



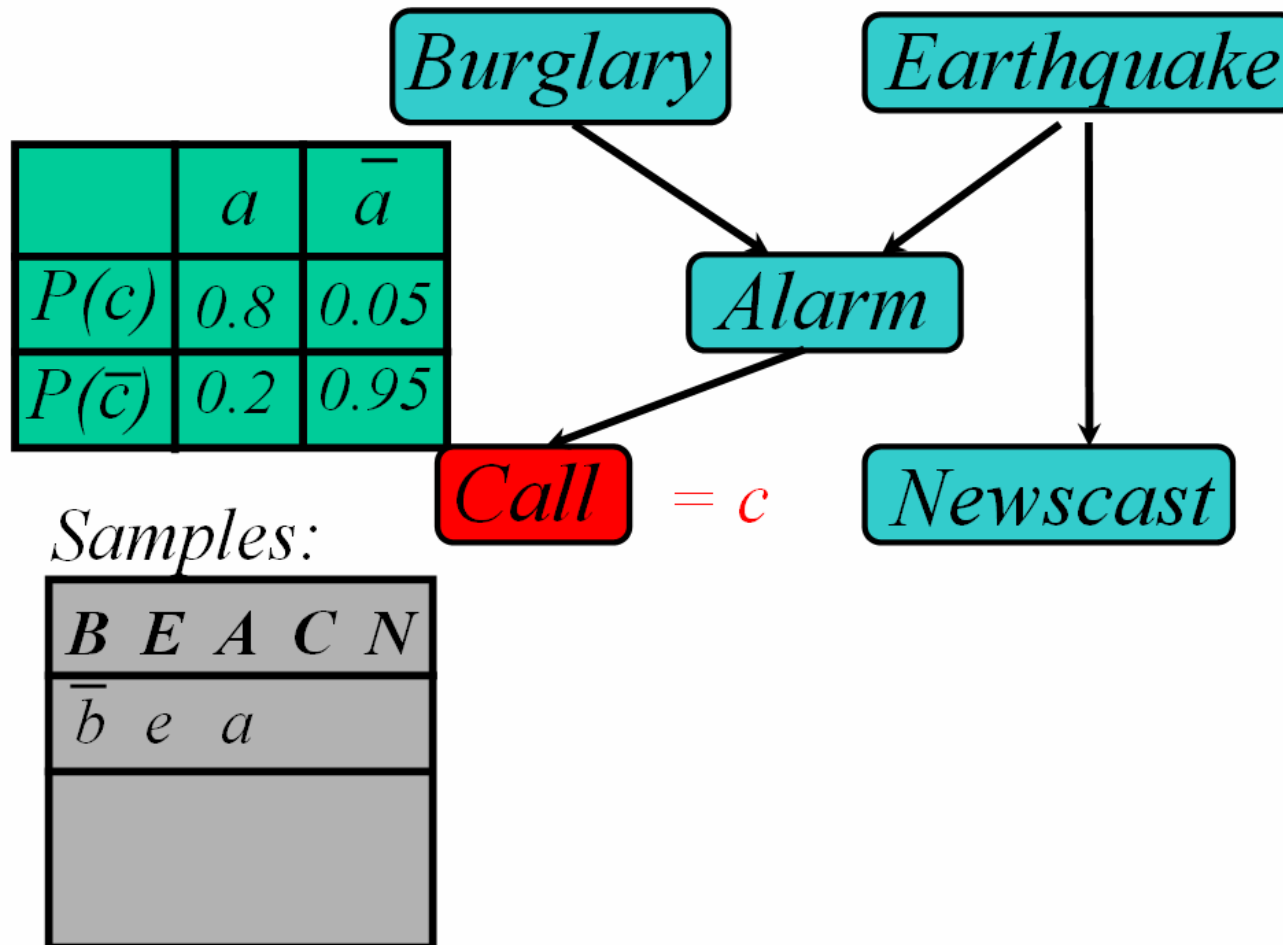
Likelihood weighting



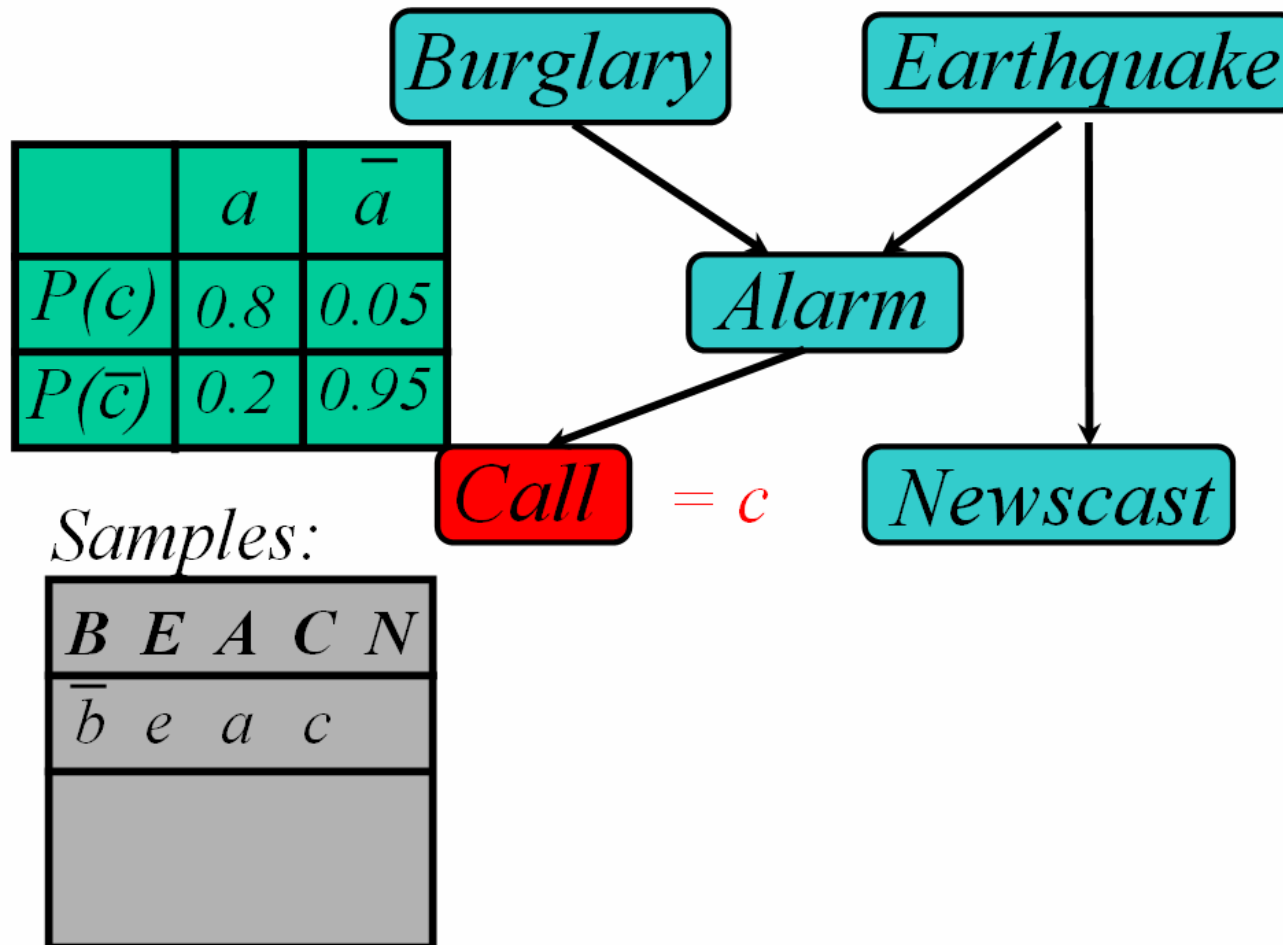
Likelihood weighting



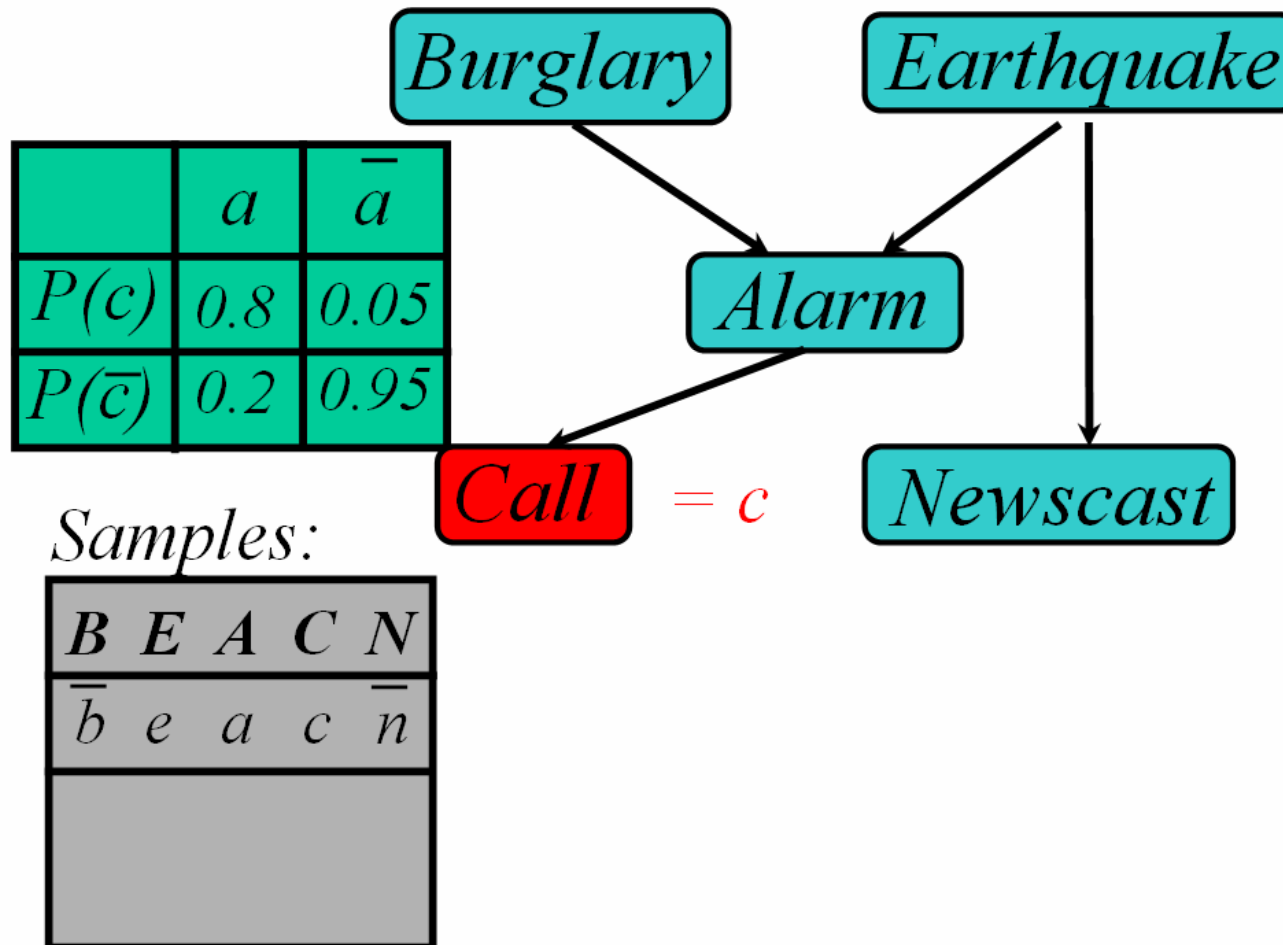
Likelihood weighting



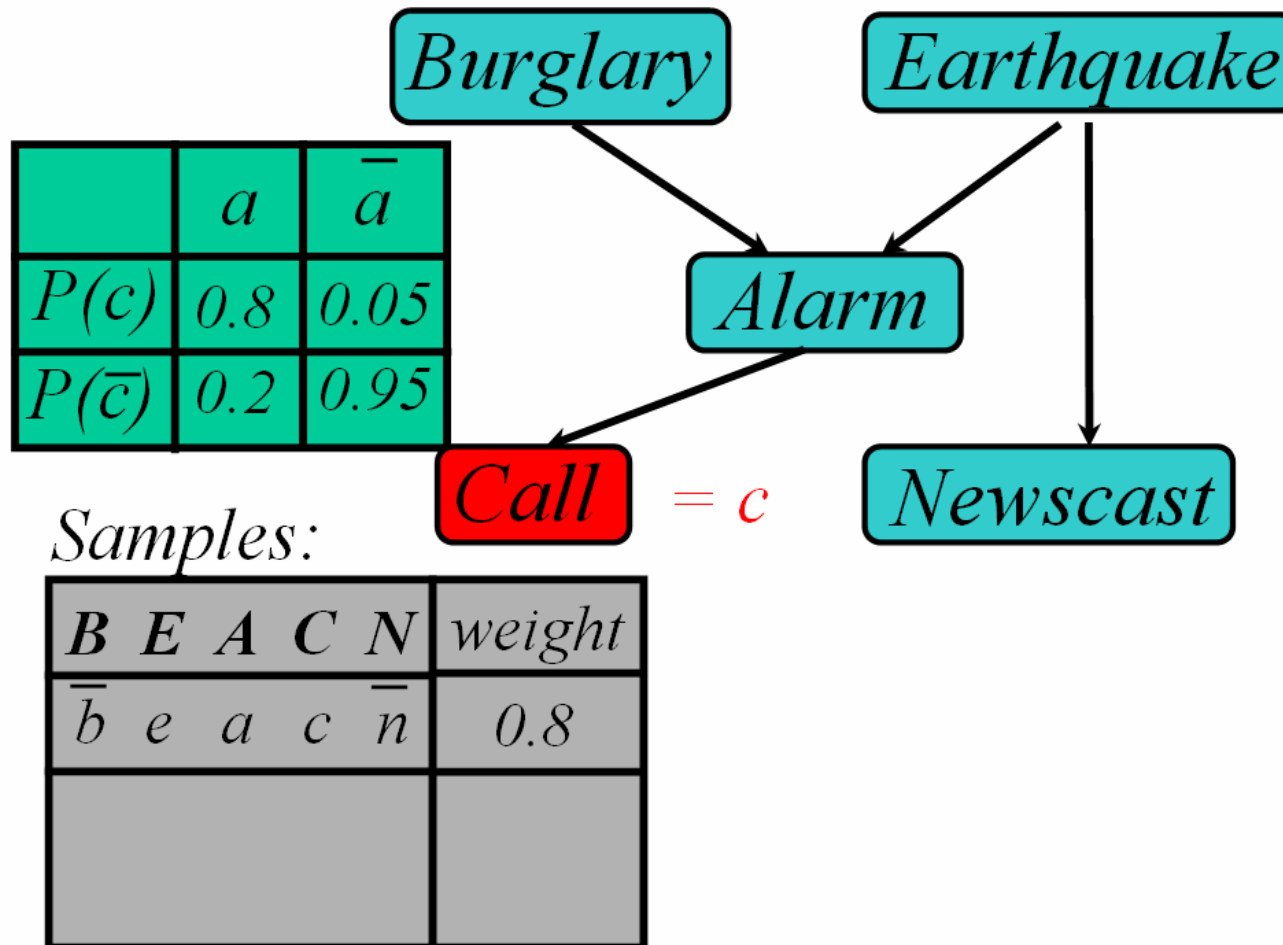
Likelihood weighting



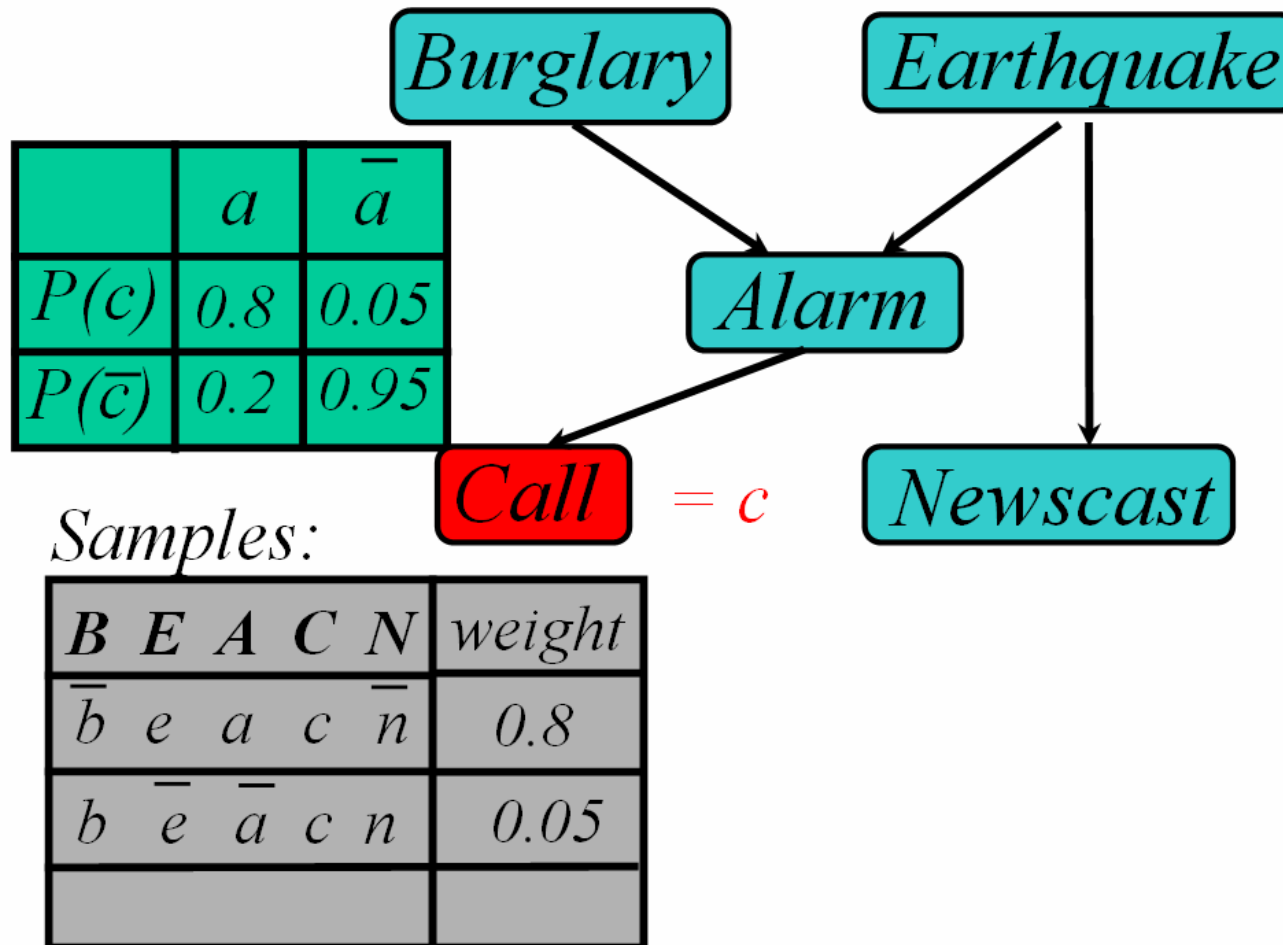
Likelihood weighting



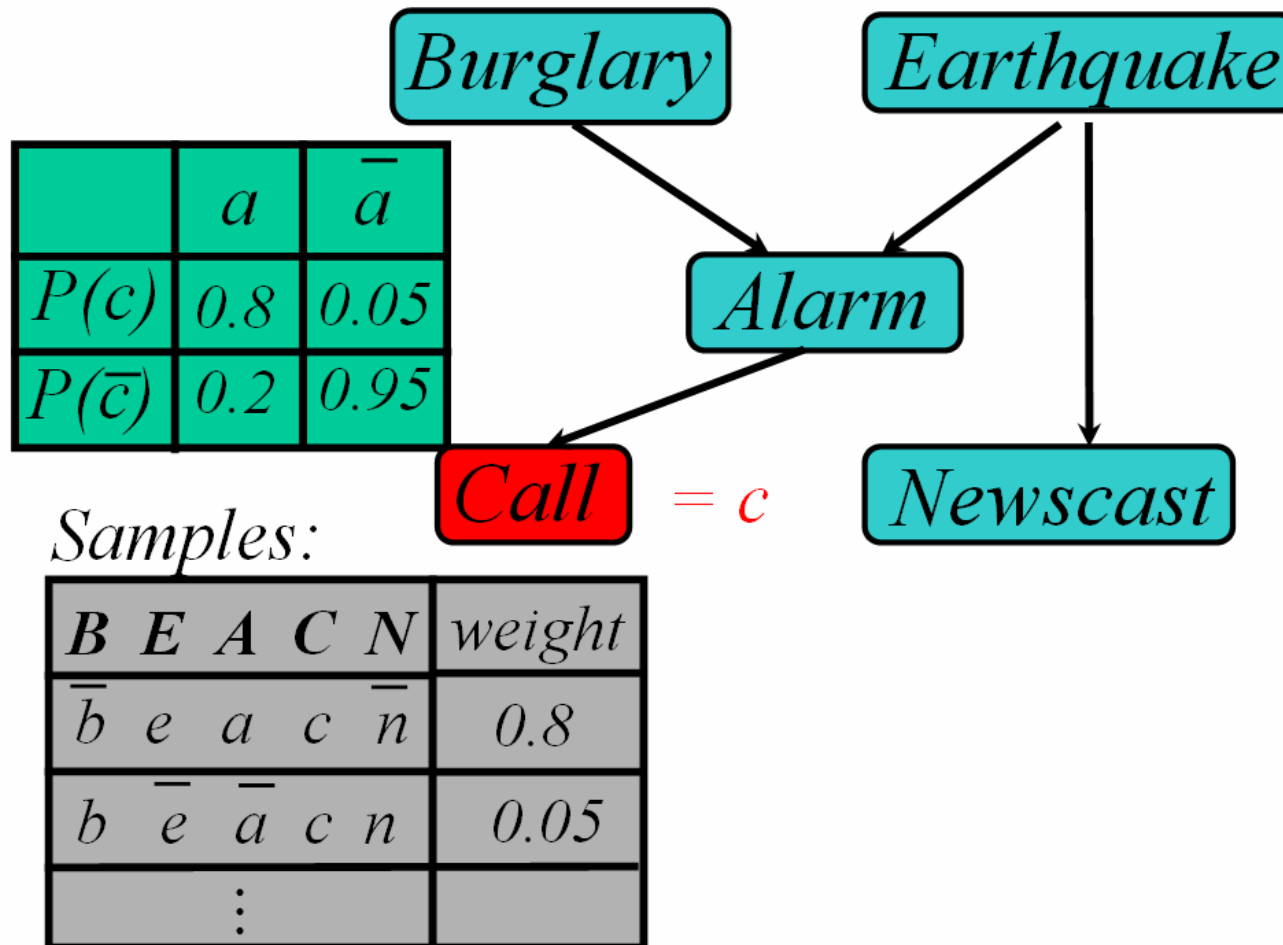
Likelihood weighting



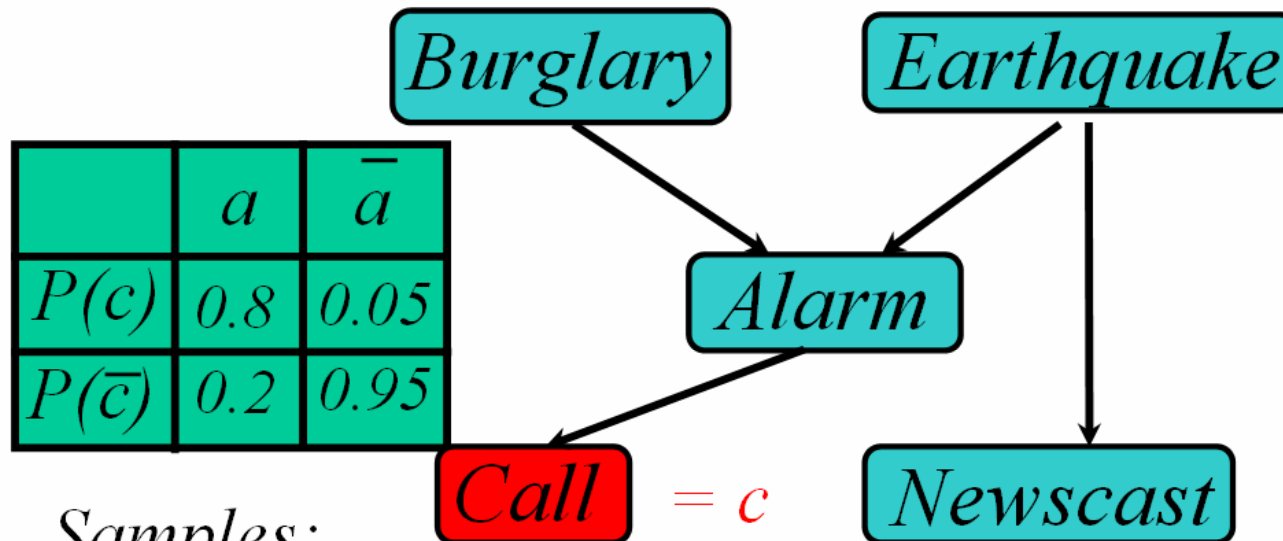
Likelihood weighting



Likelihood weighting



Likelihood weighting



Samples:

B	E	A	C	N	weight
\bar{b}	e	a	c	\bar{n}	0.8
b	\bar{e}	\bar{a}	c	n	0.05
	\vdots				

$$P(b|c) = \frac{\text{weight of samples with } B=b}{\text{total weight of samples}}$$

Markov Chain Monte Carlo



MCMC with Gibbs Sampling

Fix the values of observed variables

Set the values of all non-observed variables randomly

Perform a random walk through the space of complete variable assignments. On each move:

1. Pick a variable X
2. Calculate $\Pr(X=\text{true} \mid \text{all other variables})$
3. Set X to true with that probability

Repeat many times. Frequency with which any variable X is true is its posterior probability.

Converges to true posterior when frequencies stop changing significantly

- stable distribution, mixing

Markov Blanket Sampling

How to calculate $\Pr(X=\text{true} \mid \text{all other variables})$?

Recall: a variable is independent of all others given it's Markov Blanket

- parents
- children
- other parents of children

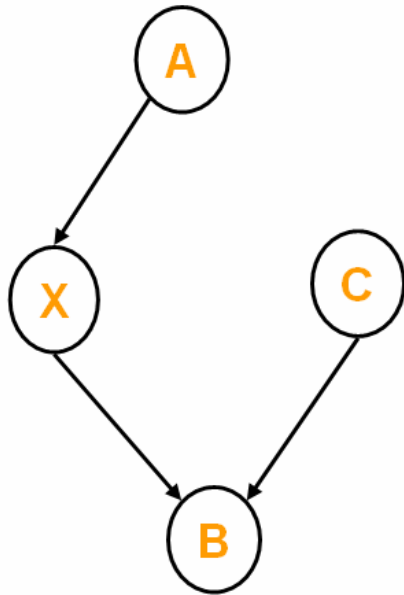
So problem becomes calculating $\Pr(X=\text{true} \mid \text{MB}(X))$

- We solve this sub-problem exactly
- Fortunately, it is easy to solve

$$P(X) = \alpha P(X \mid \text{Parents}(X)) \prod_{Y \in \text{Children}(X)} P(Y \mid \text{Parents}(Y))$$

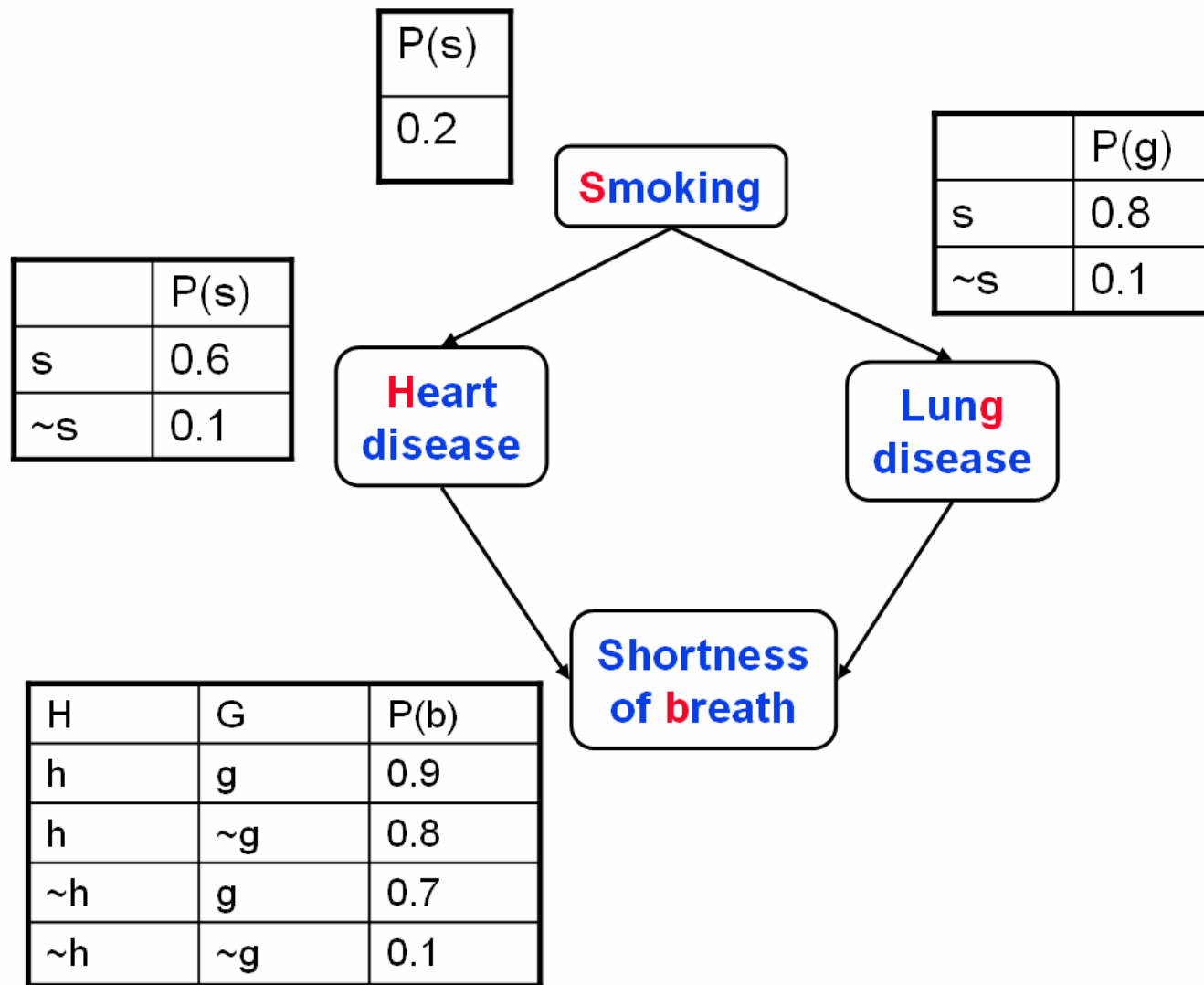
Example

$$P(X) = \alpha P(X | \text{Parents}(X)) \prod_{Y \in \text{Children}(X)} P(Y | \text{Parents}(Y))$$

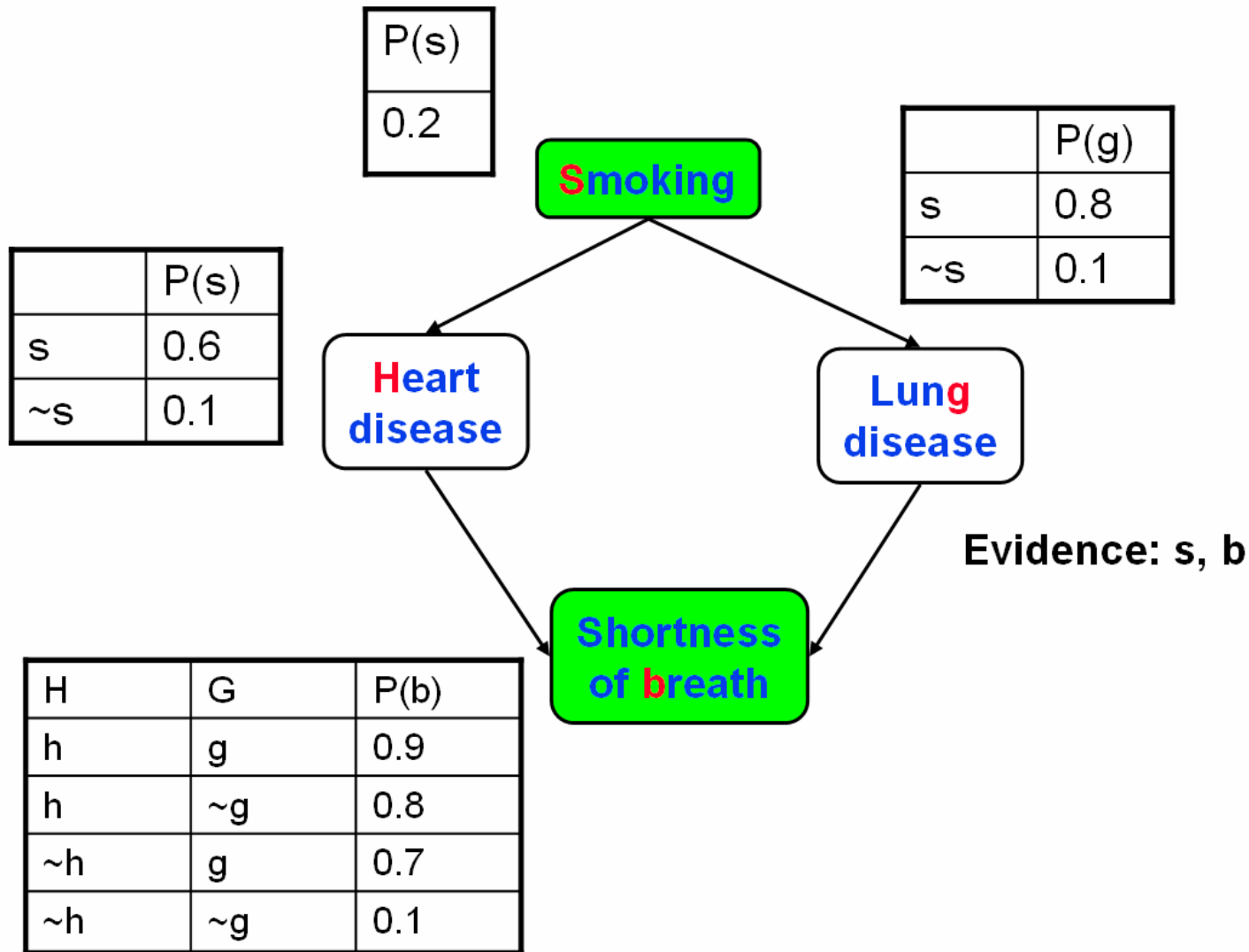


$$\begin{aligned} P(X | A, B, C) &= \frac{P(X, A, B, C)}{P(A, B, C)} \\ &= \frac{P(A)P(X | A)P(C)P(B | X, C)}{P(A, B, C)} \\ &= \left[\frac{P(A)P(C)}{P(A, B, C)} \right] P(X | A)P(B | X, C) \\ &= \alpha P(X | A)P(B | X, C) \end{aligned}$$

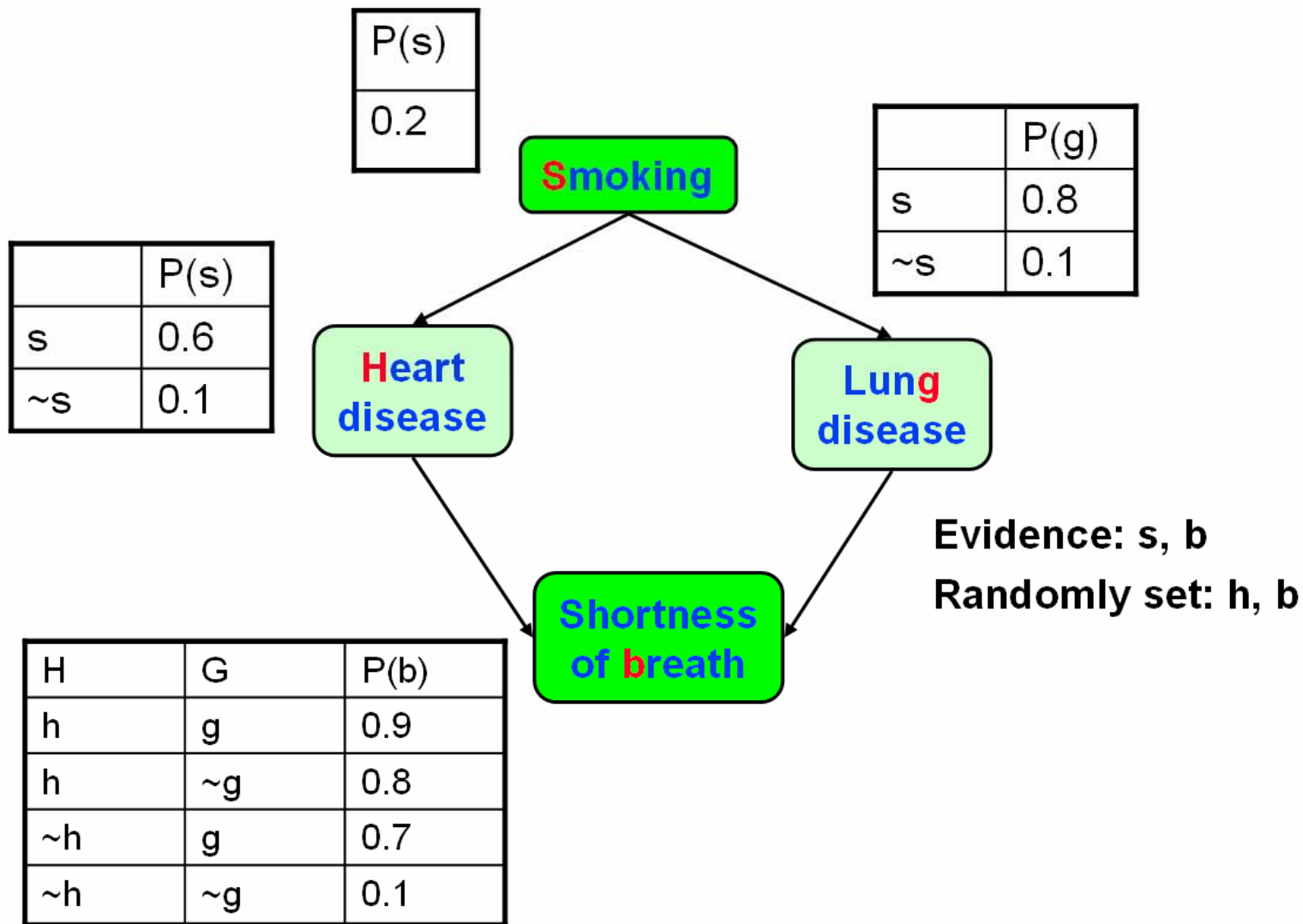
Example



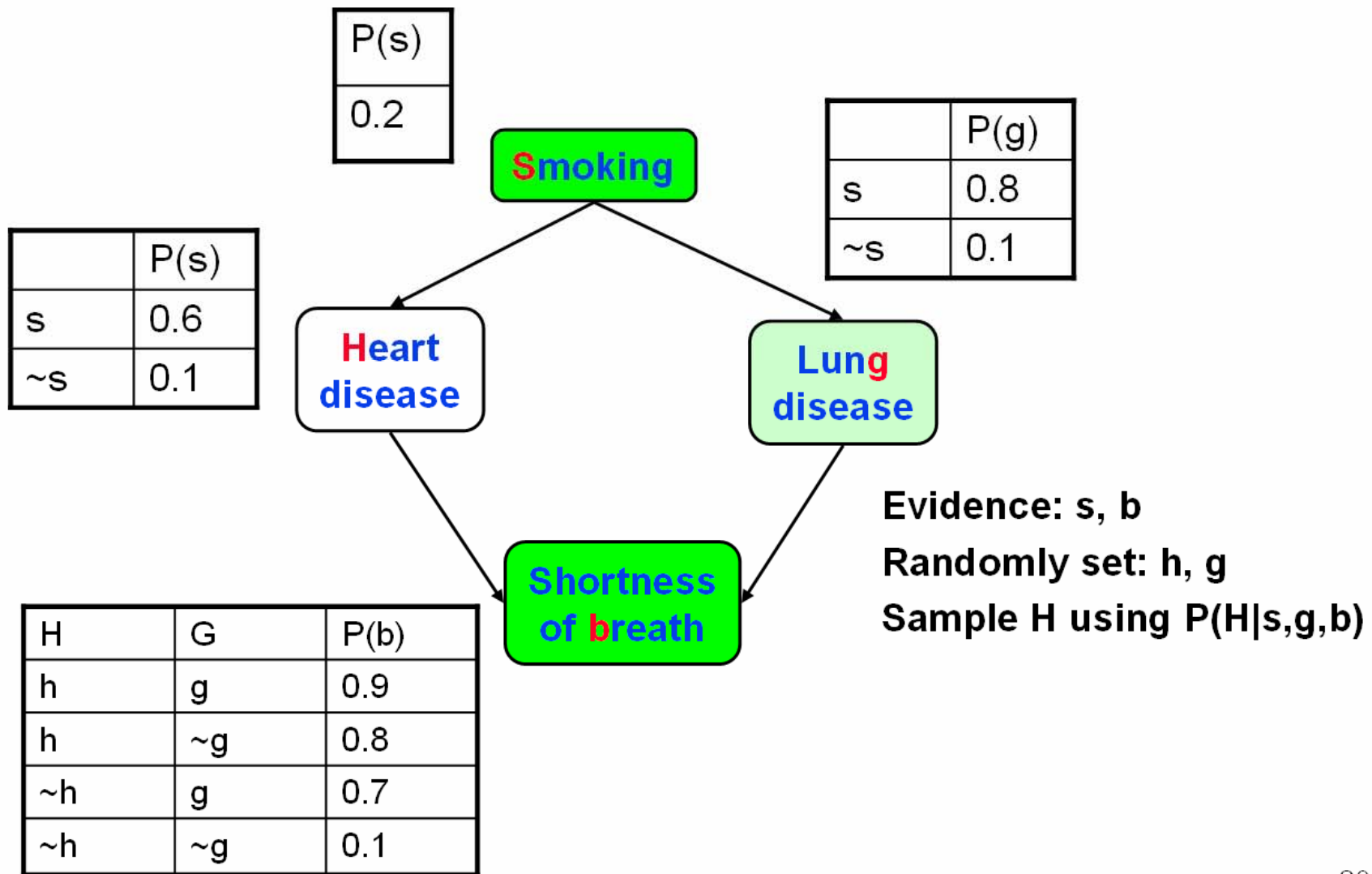
Example



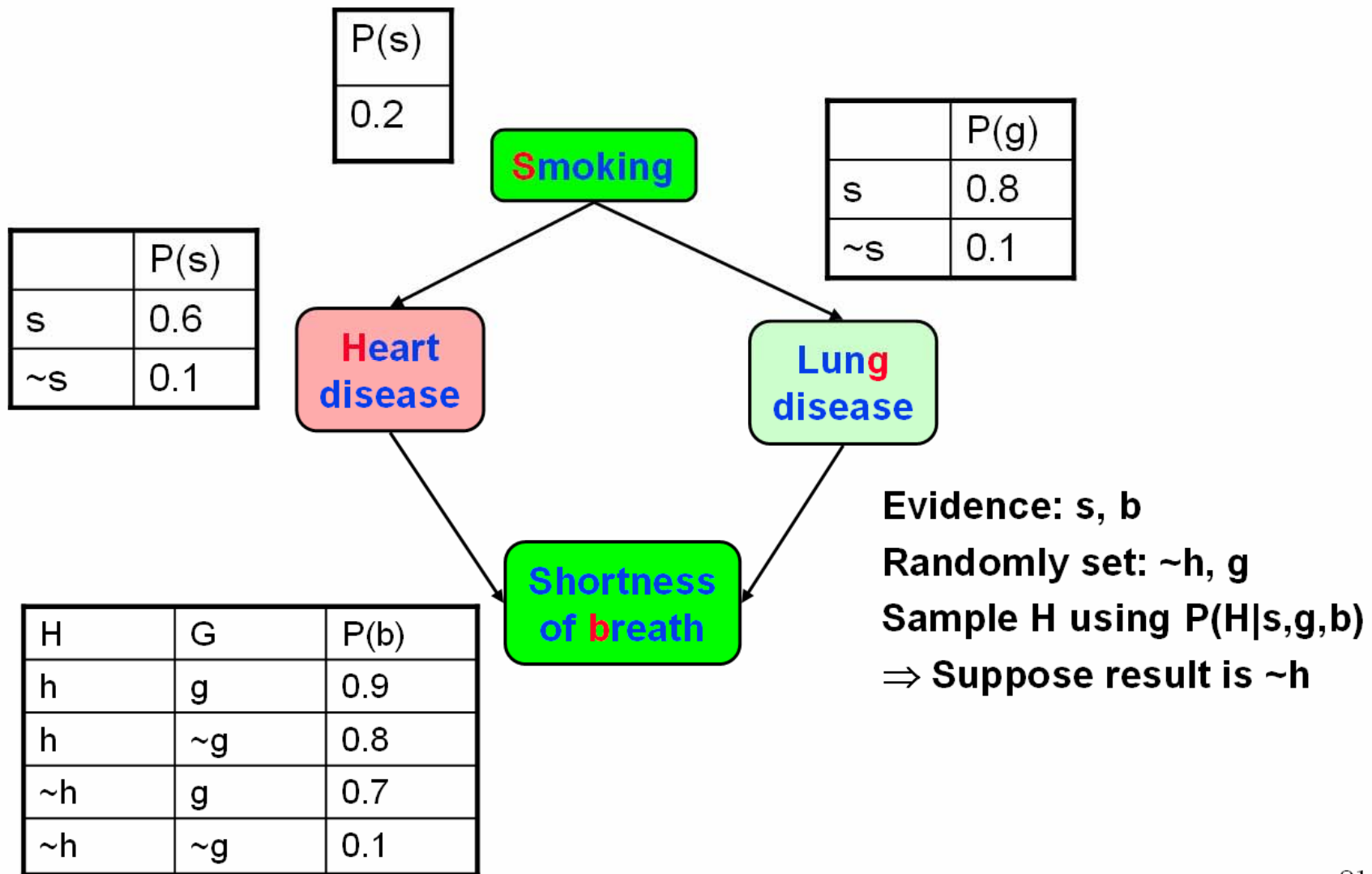
Example



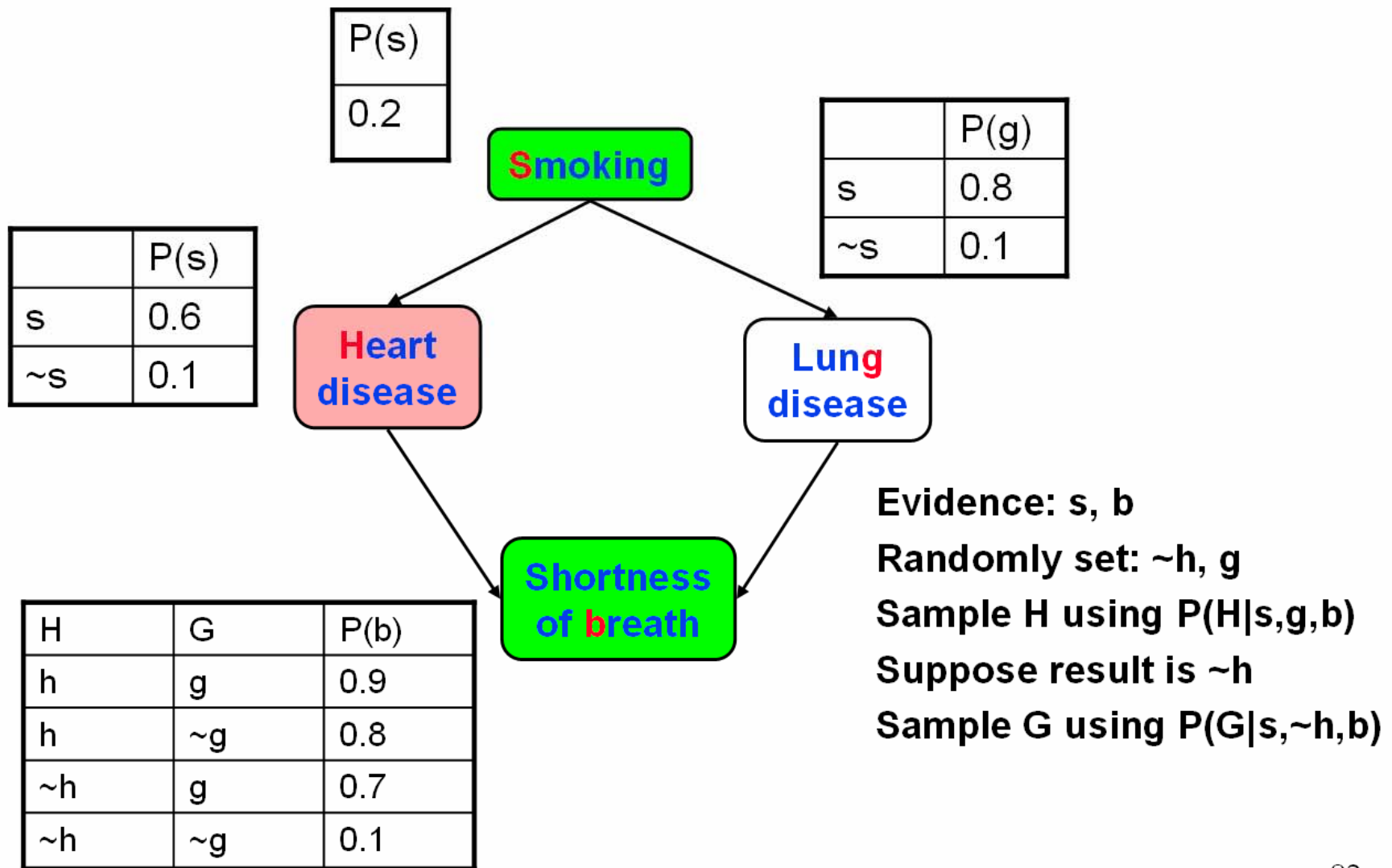
Example



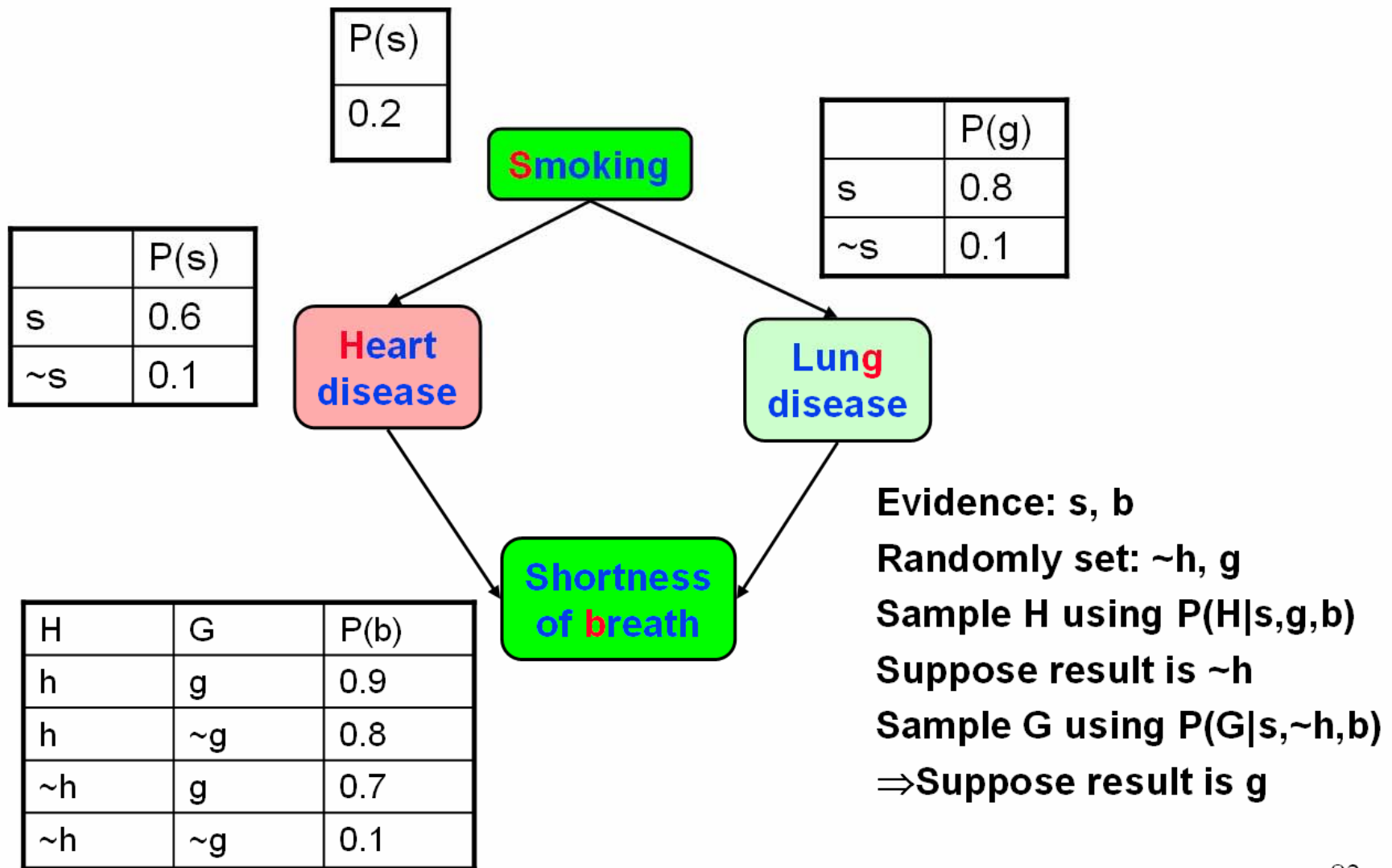
Example



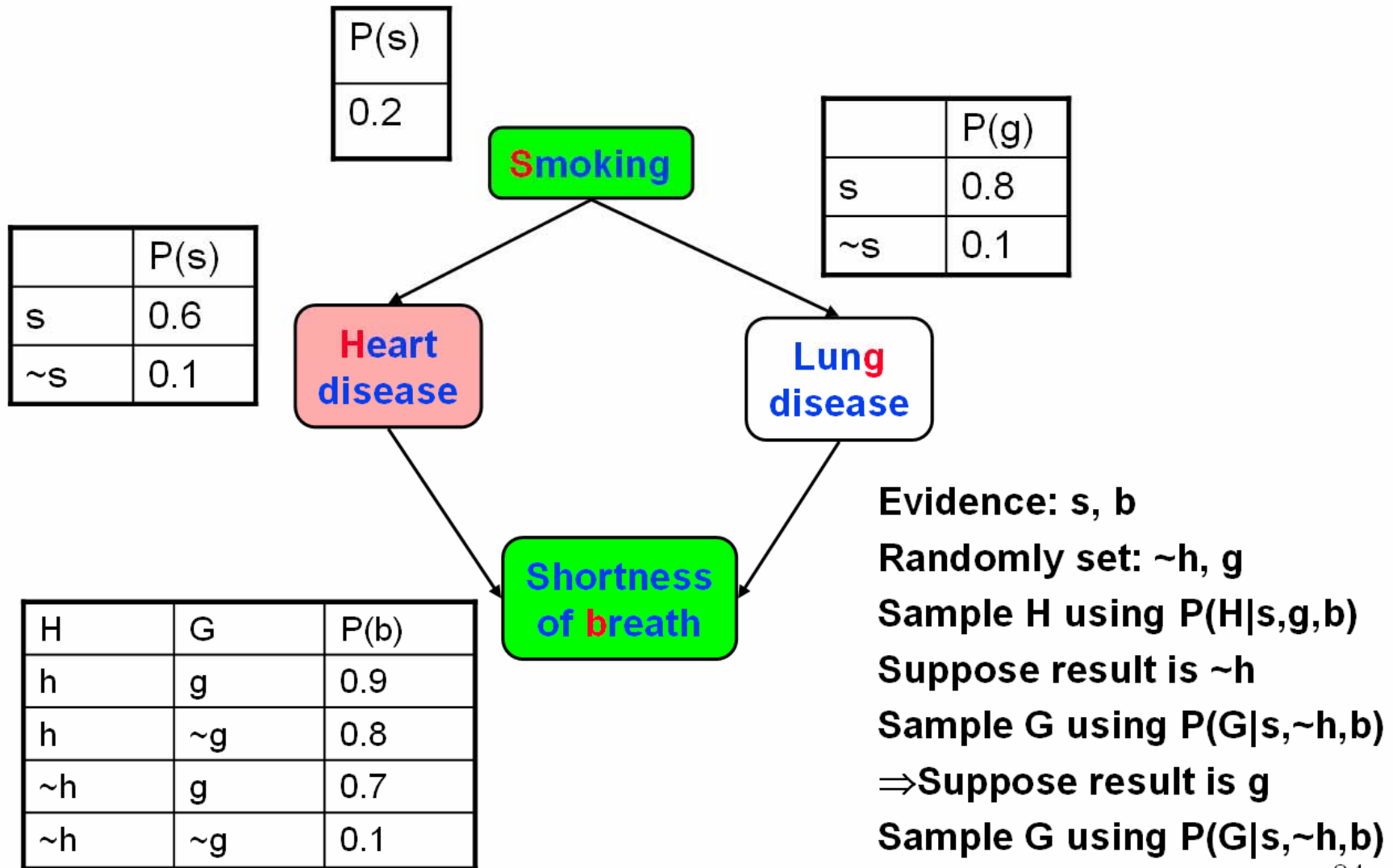
Example



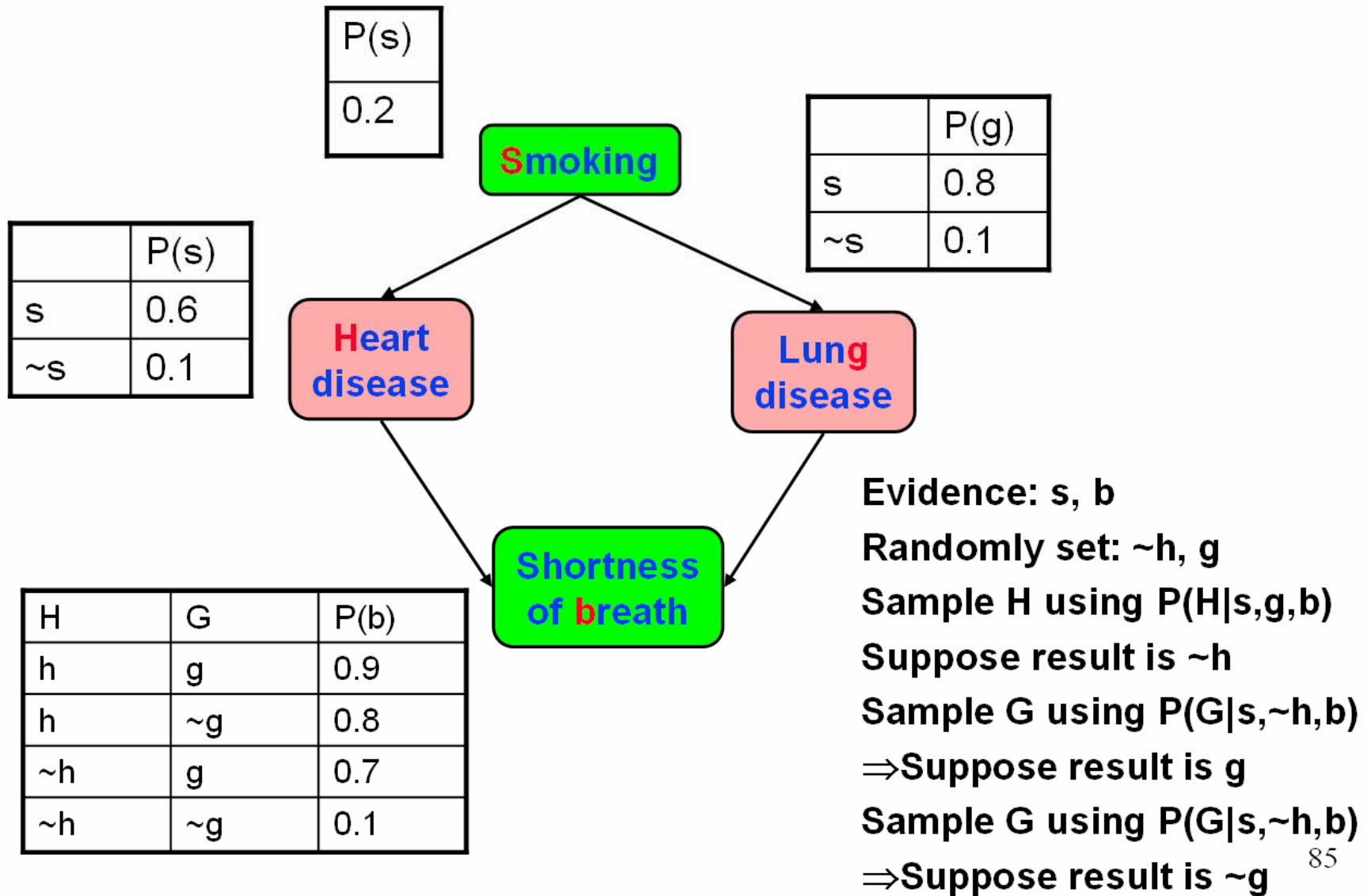
Example



Example



Example



Gibbs MCMC Summary

Advantages:

- No samples are discarded
- No problem with samples of low weight
- Can be implemented very efficiently
 - 10K samples @ second

Disadvantages:

- Can get stuck if relationship between two variables is *deterministic*
- Many variations have been devised to make MCMC more robust

Gibbs MCMC Summary

$$P(X|E) = \frac{\text{number of samples with } X=x}{\text{total number of samples}}$$

Advantages:

- No samples are discarded
- No problem with samples of low weight
- Can be implemented very efficiently
 - 10K samples @ second

Disadvantages:

- Can get stuck if relationship between two variables is *deterministic*
- Many variations have been devised to make MCMC more robust

Other approaches

- Search based techniques
 - ◆ search for high-probability instantiations
 - ◆ use instantiations to approximate probabilities
- Structural approximation
 - ◆ simplify network
 - eliminate edges, nodes
 - abstract node values
 - simplify CPTs
 - ◆ do inference in simplified network

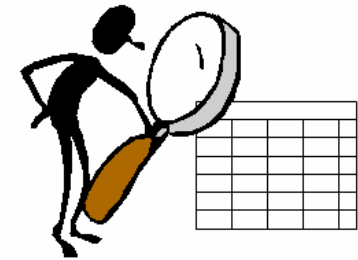
Course Contents

- Concepts in Probability
- Bayesian Networks
- Inference
- Decision making
- » Learning networks from data
- Reasoning over time
- Applications

Learning networks from data

- The learning task
- Parameter learning
 - ◆ Fully observable
 - ◆ Partially observable
- Structure learning
- Hidden variables

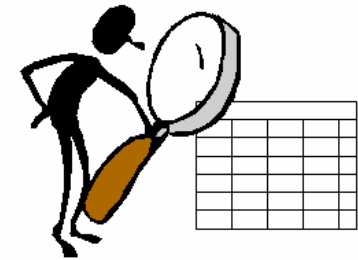
The learning task



<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}	<i>e</i>	<i>a</i>	<i>c</i>	\bar{n}
<i>b</i>	\bar{e}	\bar{a}	\bar{c}	<i>n</i>
	\vdots			

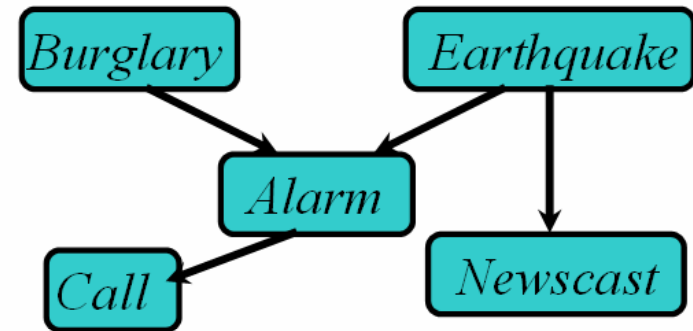
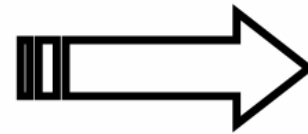
Input: training data

The learning task



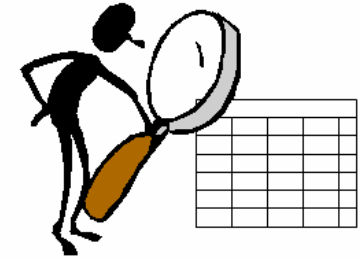
<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}	\bar{e}	\bar{a}	\bar{c}	\bar{n}
<i>b</i>	\bar{e}	\bar{a}	\bar{c}	<i>n</i>
				⋮

Input: training data

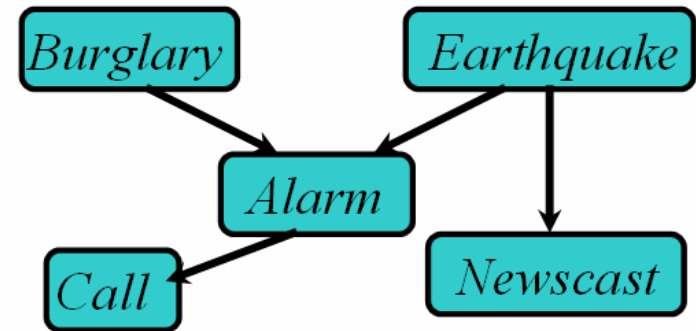
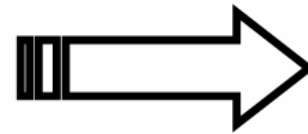


Output: BN modeling data

The learning task



<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}	\bar{e}	\bar{a}	\bar{c}	\bar{n}
<i>b</i>	\bar{e}	\bar{a}	\bar{c}	<i>n</i>
				\vdots

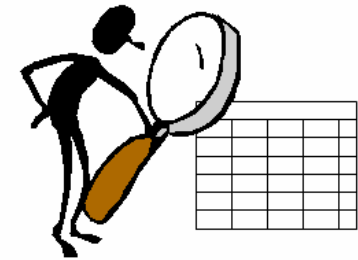


Input: training data

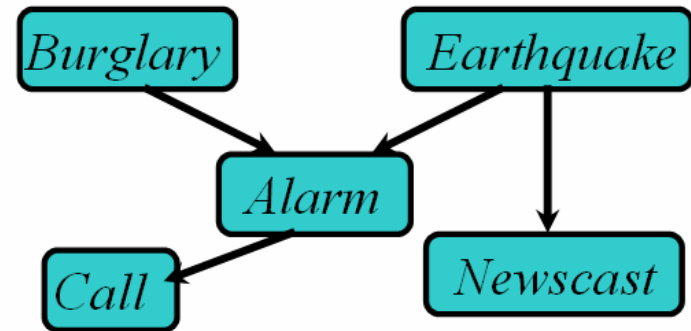
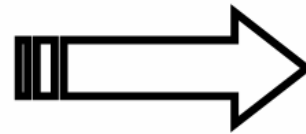
Output: BN modeling data

- Input: fully or partially observable data cases?

The learning task



<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}	\bar{e}	\bar{a}	\bar{c}	\bar{n}
<i>b</i>	\bar{e}	\bar{a}	\bar{c}	<i>n</i>
				⋮



Input: training data

Output: BN modeling data

- Input: fully or partially observable data cases?
- Output: parameters or also structure?

Parameter learning: one variable



Parameter learning: one variable



- Unfamiliar coin:

Parameter learning: one variable



- Unfamiliar coin:

- ◆ Let $\theta =$ bias of coin (long-run fraction of heads)

Parameter learning: one variable



- Unfamiliar coin:
 - ◆ Let θ = bias of coin (long-run fraction of heads)
- If θ known (given), then

Parameter learning: one variable



- Unfamiliar coin:

- ◆ Let θ = bias of coin (long-run fraction of heads)

- If θ known (given), then

- ◆ $P(X = heads \mid \theta) =$

Parameter learning: one variable



- Unfamiliar coin:
 - ◆ Let θ = bias of coin (long-run fraction of heads)
- If θ known (given), then
 - ◆ $P(X = \textit{heads} \mid \theta) = \theta$

Parameter learning: one variable



- Unfamiliar coin:
 - ◆ Let θ = bias of coin (long-run fraction of heads)
- If θ known (given), then
 - ◆ $P(X = \text{heads} \mid \theta) = \theta$
- Different coin tosses independent given θ
 - ⇒ $P(X_1, \dots, X_n \mid \theta) =$

Parameter learning: one variable



- Unfamiliar coin:
 - ◆ Let θ = bias of coin (long-run fraction of heads)
- If θ known (given), then
 - ◆ $P(X = \text{heads} \mid \theta) = \theta$
- Different coin tosses independent given θ
 - ⇒ $P(\underbrace{X_1, \dots, X_n}_{h \text{ heads, } t \text{ tails}} \mid \theta) =$

Parameter learning: one variable



- Unfamiliar coin:

- ◆ Let θ = bias of coin (long-run fraction of heads)

- If θ known (given), then

- ◆ $P(X = \text{heads} \mid \theta) = \theta$

- Different coin tosses independent given θ

- $\Rightarrow P(\underbrace{X_1, \dots, X_n}_{h \text{ heads, } t \text{ tails}} \mid \theta) = \theta^h (1-\theta)^t$

Maximum likelihood

Maximum likelihood

- Input: a set of previous coin tosses
 - ◆ $X_1, \dots, X_n = \{H, T, H, H, H, T, T, H, \dots, H\}$

Maximum likelihood

- Input: a set of previous coin tosses

- ◆ $X_1, \dots, X_n = \{ \underbrace{\text{H, T, H, H, H, T, T, H, \dots, H}}_{h \text{ heads, } t \text{ tails}} \}$

Maximum likelihood

- Input: a set of previous coin tosses

- ◆ $X_1, \dots, X_n = \{\underbrace{\text{H, T, H, H, H, T, T, H, \dots, H}}_{h \text{ heads, } t \text{ tails}}\}$

h heads, *t* tails

- Goal: estimate θ

Maximum likelihood

- Input: a set of previous coin tosses

- ◆ $X_1, \dots, X_n = \{\underbrace{\text{H, T, H, H, H, T, T, H, \dots, H}}_{h \text{ heads, } t \text{ tails}}\}$

- Goal: estimate θ

- The likelihood $P(X_1, \dots, X_n | \theta) = \theta^h (1-\theta)^t$

Maximum likelihood

- Input: a set of previous coin tosses

- ◆ $X_1, \dots, X_n = \{\underbrace{\text{H, T, H, H, H, T, T, H, \dots, H}}_{h \text{ heads, } t \text{ tails}}\}$

- Goal: estimate θ

- The likelihood $P(X_1, \dots, X_n | \theta) = \theta^h (1-\theta)^t$

- The maximum likelihood solution is:

Maximum likelihood

- Input: a set of previous coin tosses
 - ◆ $X_1, \dots, X_n = \{\underbrace{\text{H, T, H, H, H, T, T, H, \dots, H}}_{h \text{ heads, } t \text{ tails}}\}$
- Goal: estimate θ
- The likelihood $P(X_1, \dots, X_n | \theta) = \theta^h (1-\theta)^t$
- The maximum likelihood solution is:

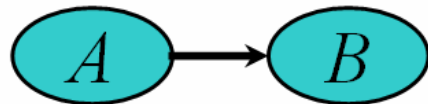
$$\theta^* = \frac{h}{h+t}$$

General parameter learning

- A multi-variable BN is composed of several independent parameters (“coins”).

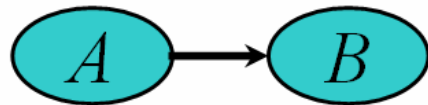
General parameter learning

- A multi-variable BN is composed of several independent parameters (“coins”).



General parameter learning

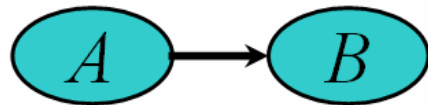
- A multi-variable BN is composed of several independent parameters (“coins”).



Three parameters:

General parameter learning

- A multi-variable BN is composed of several independent parameters (“coins”).

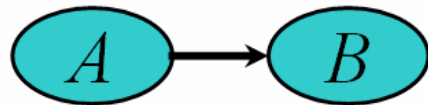


Three parameters:

$$\theta_A, \theta_{B|a}, \theta_{B|\bar{a}}$$

General parameter learning

- A multi-variable BN is composed of several independent parameters (“coins”).



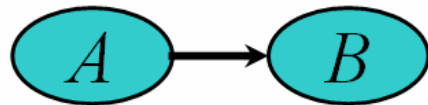
Three parameters:

$$\theta_A, \theta_{B|a}, \theta_{B|\bar{a}}$$

- Can use same techniques as one-variable case to learn each one separately

General parameter learning

- A multi-variable BN is composed of several independent parameters (“coins”).



Three parameters:

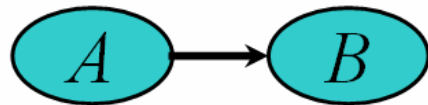
$$\theta_A, \theta_{B|a}, \theta_{B|\bar{a}}$$

- Can use same techniques as one-variable case to learn each one separately

Max likelihood estimate of $\theta_{B|\bar{a}}$ would be:

General parameter learning

- A multi-variable BN is composed of several independent parameters (“coins”).



Three parameters:

$$\theta_A, \theta_{B|a}, \theta_{B|\bar{a}}$$

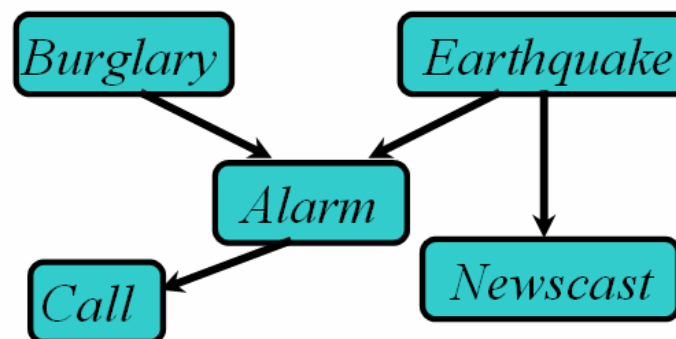
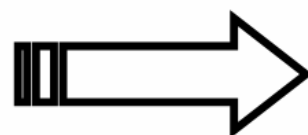
- Can use same techniques as one-variable case to learn each one separately

Max likelihood estimate of $\theta_{B|\bar{a}}$ would be:

$$\theta_{B|\bar{a}}^* = \frac{\#data\ cases\ with\ b, \bar{a}}{\#data\ cases\ with\ \bar{a}}$$

Partially observable data

<i>B</i>	<i>E</i>	<i>A</i>	<i>C</i>	<i>N</i>
\bar{b}	?	<i>a</i>	<i>c</i>	?
<i>b</i>	?	\bar{a}	?	<i>n</i>
		⋮		



- Fill in missing data with “expected” value
 - ◆ expected = distribution over possible values
 - ◆ use “best guess” BN to estimate distribution

Intuition

- In fully observable case:

Intuition

- In fully observable case:

$$\theta_{n|e}^* = \frac{\#data\ cases\ with\ n,\ e}{\#data\ cases\ with\ e}$$

Intuition

- In fully observable case:

$$\theta_{n|e}^* = \frac{\#data\ cases\ with\ n,\ e}{\#data\ cases\ with\ e} = \frac{\sum_j I(n, e | d_j)}{\sum_j I(e | d_j)}$$

$$I(e | d_j) = \begin{cases} 1 & \text{if } E=e \text{ in data case } d_j \\ 0 & \text{otherwise} \end{cases}$$

Intuition

- In fully observable case:

$$\theta_{n|e}^* = \frac{\#data\ cases\ with\ n,\ e}{\#data\ cases\ with\ e} = \frac{\sum_j I(n, e | d_j)}{\sum_j I(e | d_j)}$$

$$I(e | d_j) = \begin{cases} 1 & \text{if } E=e \text{ in data case } d_j \\ 0 & \text{otherwise} \end{cases}$$

- In partially observable case I is unknown.

Intuition

- In fully observable case:

$$\theta_{n|e}^* = \frac{\#data\ cases\ with\ n,\ e}{\#data\ cases\ with\ e} = \frac{\sum_j I(n, e | d_j)}{\sum_j I(e | d_j)}$$

$$I(e | d_j) = \begin{cases} 1 & \text{if } E=e \text{ in data case } d_j \\ 0 & \text{otherwise} \end{cases}$$

- In partially observable case I is unknown.

Best estimate for I is:

Intuition

- In fully observable case:

$$\theta_{n|e}^* = \frac{\#data\ cases\ with\ n,\ e}{\#data\ cases\ with\ e} = \frac{\sum_j I(n, e | d_j)}{\sum_j I(e | d_j)}$$

$$I(e | d_j) = \begin{cases} 1 & \text{if } E=e \text{ in data case } d_j \\ 0 & \text{otherwise} \end{cases}$$

- In partially observable case I is unknown.

Best estimate for I is: $\hat{I}(n, e | d_j) = P_{\theta^*}(n, e | d_j)$

Intuition

- In fully observable case:

$$\theta_{n|e}^* = \frac{\#data\ cases\ with\ n,\ e}{\#data\ cases\ with\ e} = \frac{\sum_j I(n, e | d_j)}{\sum_j I(e | d_j)}$$

$$I(e | d_j) = \begin{cases} 1 & \text{if } E=e \text{ in data case } d_j \\ 0 & \text{otherwise} \end{cases}$$

- In partially observable case I is unknown.

Best estimate for I is: $\hat{I}(n, e | d_j) = P_{\theta^*}(n, e | d_j)$

Problem: θ^* unknown.

Expectation Maximization (EM)

Expectation Maximization (EM)

- Expectation (E) step
 - ◆ Use current parameters θ to estimate filled in data.

Expectation Maximization (EM)

- Expectation (E) step
 - ◆ Use current parameters θ to estimate filled in data.

$$\hat{I}(n, e | d_j) = P_{\theta} (n, e | d_j)$$

Expectation Maximization (EM)

- Expectation (E) step
 - ◆ Use current parameters θ to estimate filled in data.

$$\hat{I}(n, e | d_j) = P_{\theta} (n, e | d_j)$$

- Maximization (M) step
 - ◆ Use filled in data to do max likelihood estimation

Expectation Maximization (EM)

- Expectation (E) step
 - ◆ Use current parameters θ to estimate filled in data.

$$\hat{I}(n, e | d_j) = P_{\theta}(n, e | d_j)$$

- Maximization (M) step
 - ◆ Use filled in data to do max likelihood estimation

$$\tilde{\theta}_{n|e} = \frac{\sum_j \hat{I}(n, e | d_j)}{\sum_j \hat{I}(e | d_j)}$$

Expectation Maximization (EM)

- Expectation (E) step
 - ◆ Use current parameters θ to estimate filled in data.

$$\hat{I}(n, e | d_j) = P_{\theta}(n, e | d_j)$$

- Maximization (M) step
 - ◆ Use filled in data to do max likelihood estimation

$$\tilde{\theta}_{n|e} = \frac{\sum_j \hat{I}(n, e | d_j)}{\sum_j \hat{I}(e | d_j)}$$

- Set: $\theta := \tilde{\theta}$

Expectation Maximization (EM)

Repeat :

- Expectation (E) step
 - ◆ Use current parameters θ to estimate filled in data.

$$\hat{I}(n, e | d_j) = P_{\theta} (n, e | d_j)$$

- Maximization (M) step
 - ◆ Use filled in data to do max likelihood estimation

$$\tilde{\theta}_{n|e} = \frac{\sum_j \hat{I}(n, e | d_j)}{\sum_j \hat{I}(e | d_j)}$$

- Set: $\theta := \tilde{\theta}$

until convergence.

Structure learning

Goal:

find “good” BN structure (relative to data)

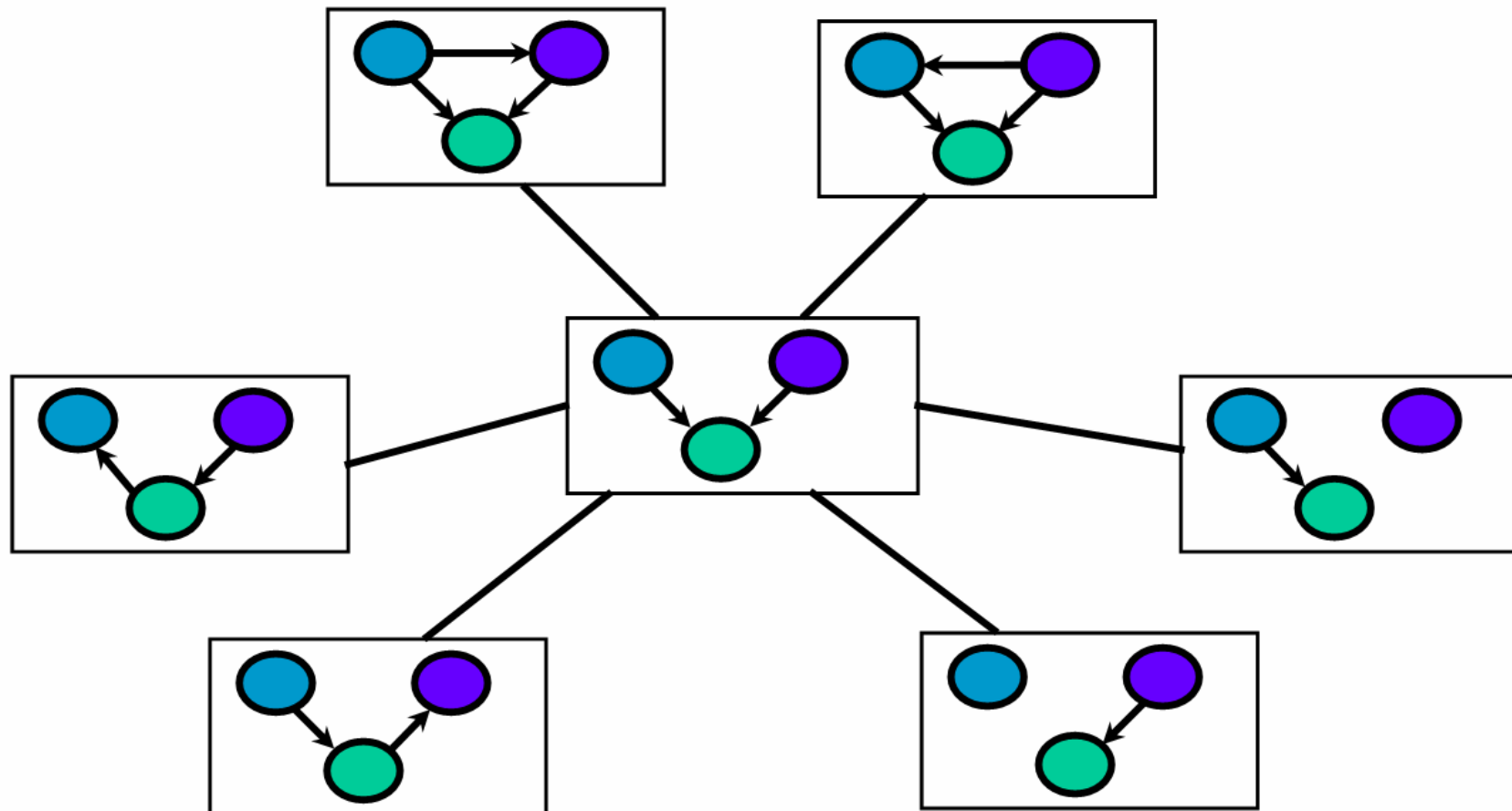
Solution:

do heuristic search over space of network structures.

Search space

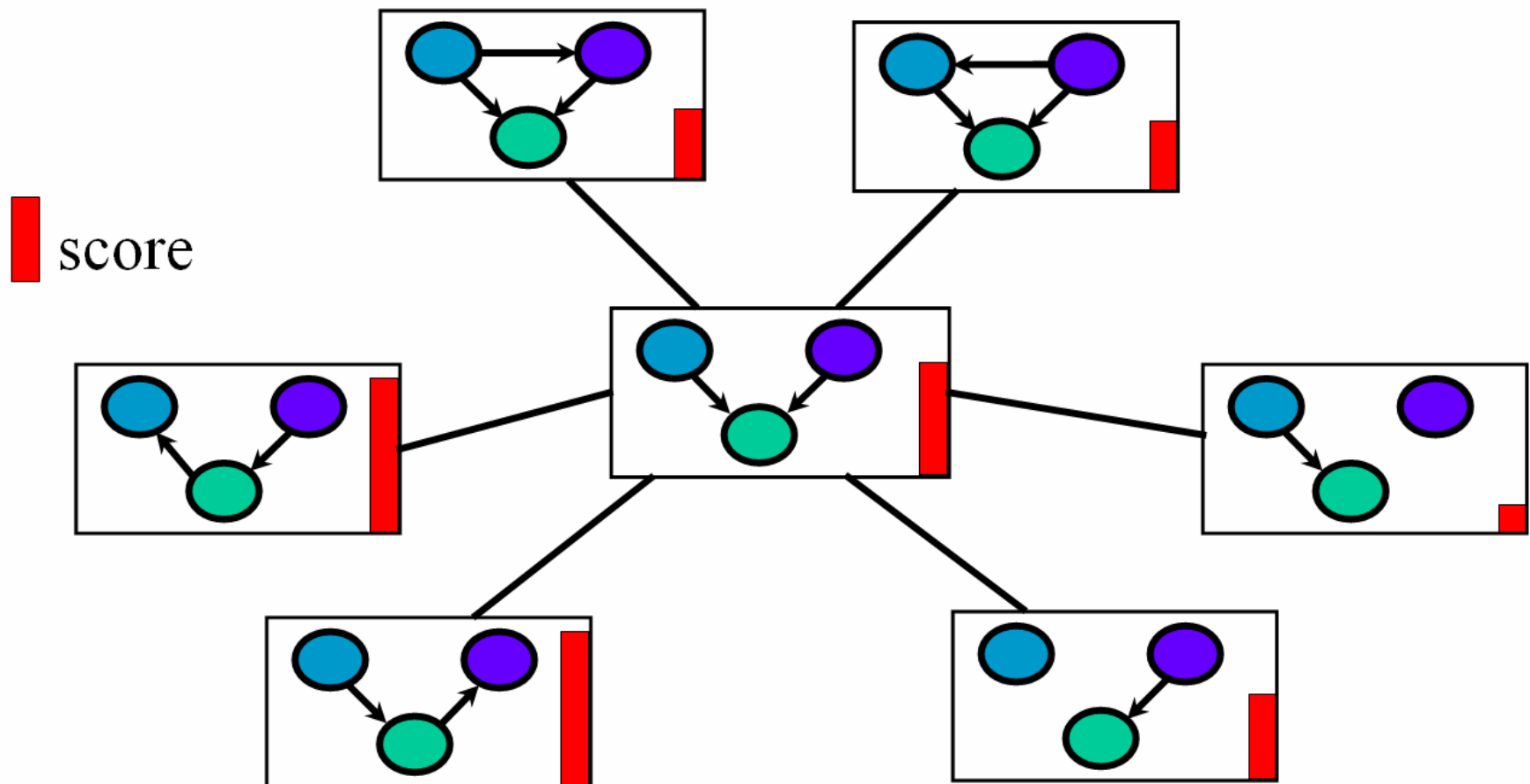
Space = network structures

Operators = add/reverse/delete edges



Heuristic search

Use scoring function to do heuristic search (any algorithm).
Greedy hill-climbing with randomness works pretty well.



Scoring

Scoring

- Fill in parameters using previous techniques & score completed networks.


Scoring

- Fill in parameters using previous techniques & score completed networks.
- One possibility for score:

Scoring


- Fill in parameters using previous techniques & score completed networks.
- One possibility for score:
likelihood function: $Score(B) = P(data | B)$

Scoring

- Fill in parameters using previous techniques & score completed networks.
- One possibility for score:
likelihood function: $Score(B) = P(data | B)$ 

Scoring

- Fill in parameters using previous techniques & score completed networks.
- One possibility for score:


likelihood function: $Score(B) = P(data | B)$ 

Example: X, Y independent coin tosses

typical $data = (27\ h-h, 22\ h-t, 25\ t-h, 26\ t-t)$

Scoring

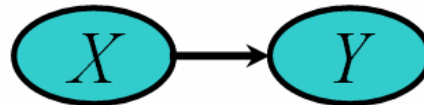
- Fill in parameters using previous techniques & score completed networks.
- One possibility for score:

likelihood function: $Score(B) = P(data | B)$ 

Example: X, Y independent coin tosses


typical $data = (27\ h-h, 22\ h-t, 25\ t-h, 26\ t-t)$

Maximum likelihood network structure:



Scoring

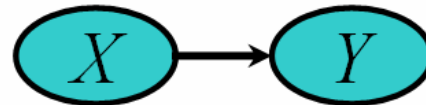
- Fill in parameters using previous techniques & score completed networks.
- One possibility for score:

likelihood function: $Score(B) = P(data | B)$ 

Example: X, Y independent coin tosses

typical $data = (27 h-h, 22 h-t, 25 t-h, 26 t-t)$

Maximum likelihood network structure:



Max. likelihood network typically fully connected

Scoring

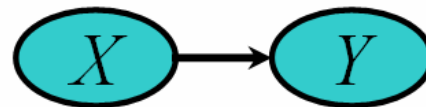
- Fill in parameters using previous techniques & score completed networks.
- One possibility for score:

likelihood function: $Score(B) = P(data | B)$ 🖐

Example: X, Y independent coin tosses

typical $data = (27\ h-h, 22\ h-t, 25\ t-h, 26\ t-t)$

Maximum likelihood network structure:



Max. likelihood network typically fully connected

This is not surprising: maximum likelihood always overfits...

Better scoring functions

Better scoring functions

- MDL formulation: balance fit to data and model complexity (# of parameters)

Better scoring functions

- MDL formulation: balance fit to data and model complexity (# of parameters)

$$\text{Score}(B) = P(\text{data} \mid B) - \text{model complexity}$$

Better scoring functions

- MDL formulation: balance fit to data and model complexity (# of parameters)

$$\text{Score}(B) = P(\text{data} \mid B) - \text{model complexity}$$

- Full Bayesian formulation

Better scoring functions

- MDL formulation: balance fit to data and model complexity (# of parameters)

$$\text{Score}(B) = P(\text{data} \mid B) - \text{model complexity}$$

- Full Bayesian formulation
 - ◆ prior on network structures & parameters

Better scoring functions

- MDL formulation: balance fit to data and model complexity (# of parameters)

$$\text{Score}(B) = P(\text{data} \mid B) - \text{model complexity}$$

- Full Bayesian formulation
 - ◆ prior on network structures & parameters
 - ◆ more parameters \Rightarrow higher dimensional space

Better scoring functions

- MDL formulation: balance fit to data and model complexity (# of parameters)

$$\text{Score}(B) = P(\text{data} \mid B) - \text{model complexity}$$

- Full Bayesian formulation
 - ◆ prior on network structures & parameters
 - ◆ more parameters \Rightarrow higher dimensional space
 - ◆ get balance effect as a byproduct*

Better scoring functions

- MDL formulation: balance fit to data and model complexity (# of parameters)

$$\text{Score}(B) = P(\text{data} \mid B) - \text{model complexity}$$

- Full Bayesian formulation
 - ◆ prior on network structures & parameters
 - ◆ more parameters \Rightarrow higher dimensional space
 - ◆ get balance effect as a byproduct*

* with Dirichlet parameter prior, MDL is an approximation to full Bayesian score.

Course Contents

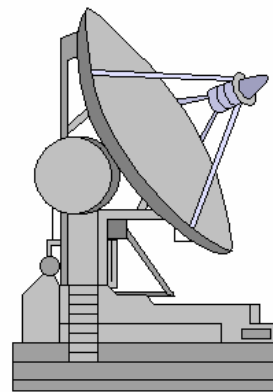
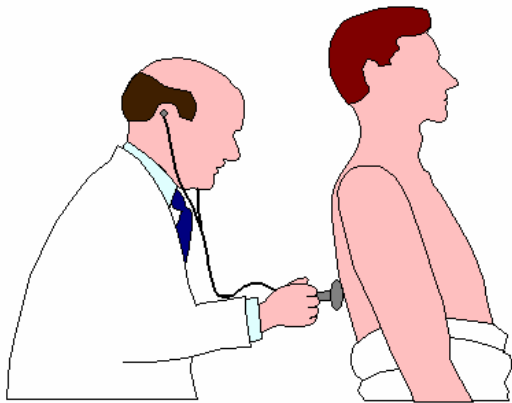
- Concepts in Probability
- Bayesian Networks
- Inference
- Decision making
- Learning networks from data
- Reasoning over time
- » Applications

Applications

- Medical expert systems
 - ◆ Pathfinder
 - ◆ Parenting MSN
- Fault diagnosis
 - ◆ Ricoh FIXIT
 - ◆ Decision-theoretic troubleshooting
- Vista
- Collaborative filtering

Why use Bayesian Networks?

- Explicit management of uncertainty/tradeoffs
- Modularity implies maintainability
- Better, flexible, and robust recommendation strategies



Pathfinder

- Pathfinder is one of the first BN systems.
- It performs diagnosis of lymph-node diseases.
- It deals with over 60 diseases and 100 findings.
- Commercialized by Intellipath and Chapman Hall publishing and applied to about 20 tissue types.

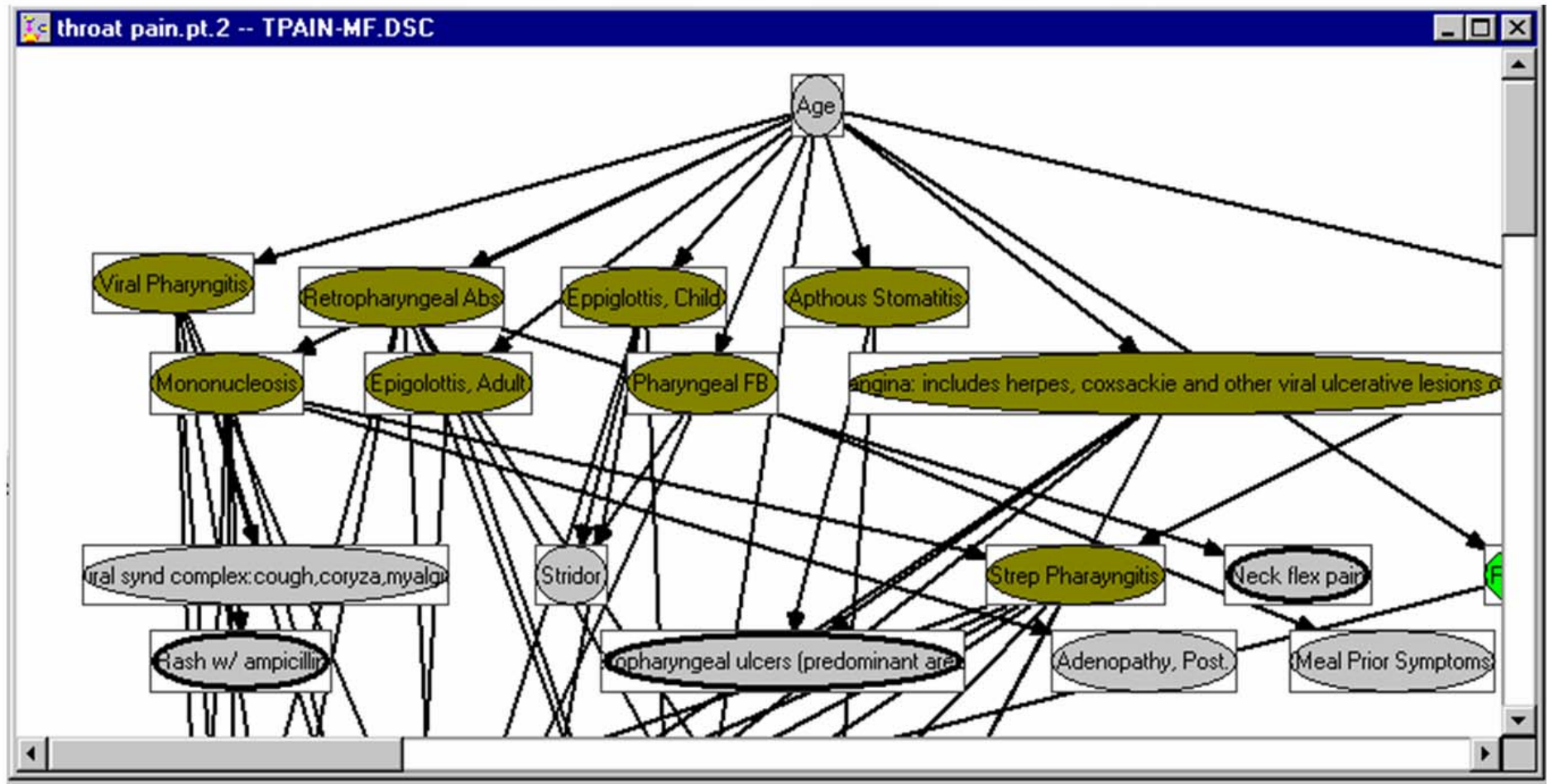
On Parenting: Selecting problem

- Diagnostic indexing for Home Health site on Microsoft Network
- Enter symptoms for pediatric complaints
- Recommends multimedia content



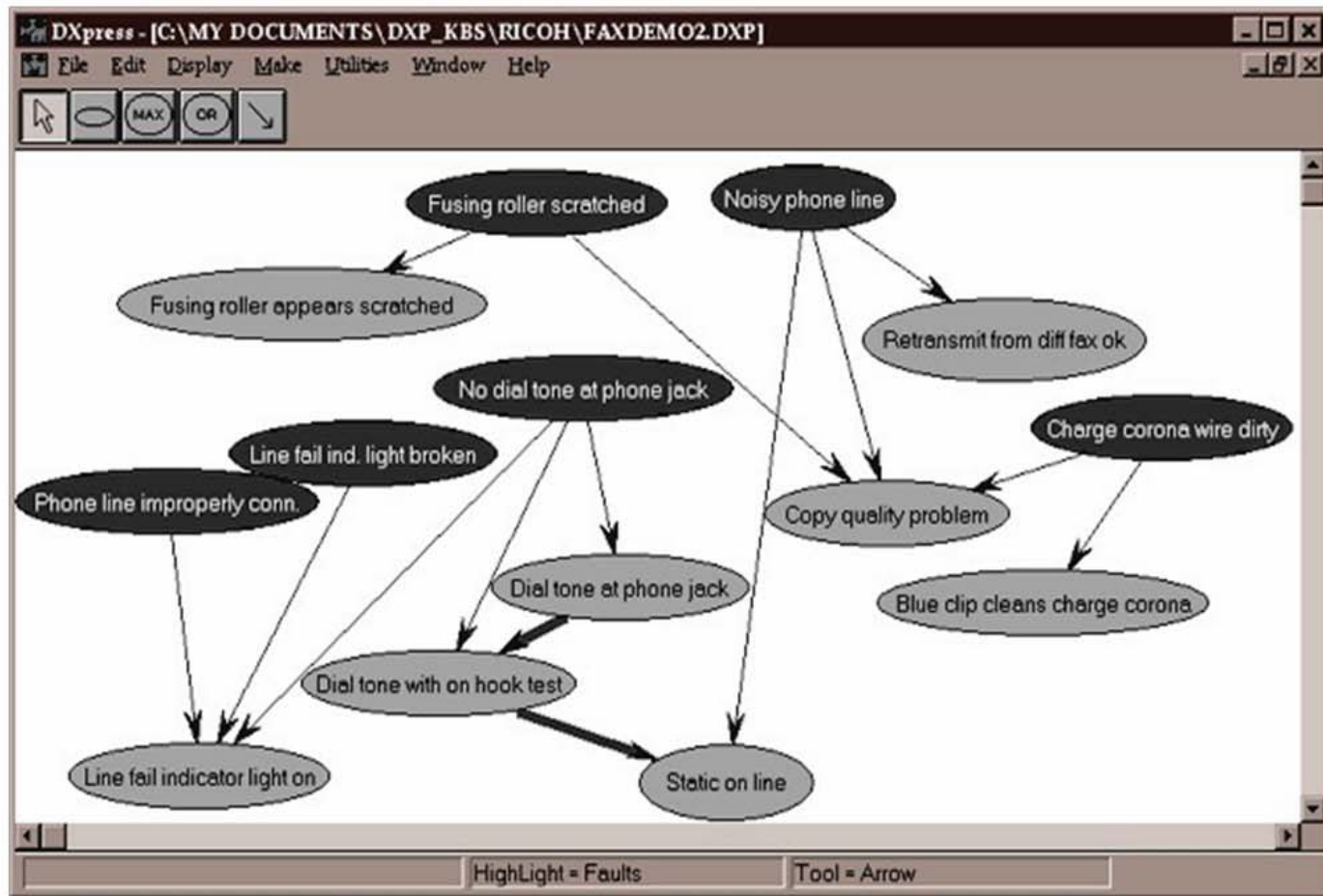
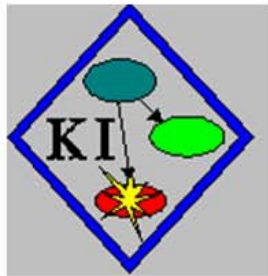
On Parenting : MSN

Original Multiple Fault Model



RICOH Fixit

- Diagnostics and information retrieval



What is Collaborative Filtering?

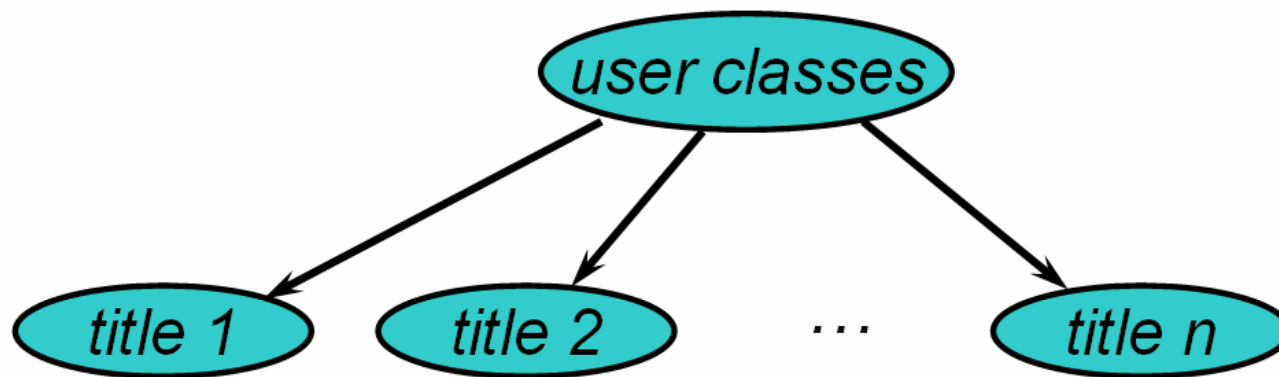
- A way to find cool websites, news stories, music artists etc
- Uses data on the preferences of many users, not descriptions of the content.
- [Firefly](#), [Net Perceptions](#) (GroupLens), and others offer this technology.

Bayesian Clustering for Collaborative Filtering

- Probabilistic summary of the data
- Reduces the number of parameters to represent a set of preferences
- Provides insight into usage patterns.
- Inference:

$$P(\text{Like title } i \mid \text{Like title } j, \text{ Like title } k)$$

Applying Bayesian clustering



	<i>class1</i>	<i>class2</i>	...
<i>title1</i>	$p(\text{like})=0.2$	$p(\text{like})=0.8$	
<i>title2</i>	$p(\text{like})=0.7$	$p(\text{like})=0.1$	
<i>title3</i>	$p(\text{like})=0.99$	$p(\text{like})=0.01$	
	...		

MSNBC Story clusters

Readers of commerce and technology stories (36%):

- E-mail delivery isn't exactly guaranteed
- Should you buy a DVD player?
- Price low, demand high for Nintendo

MSNBC Story clusters

Readers of commerce and technology stories (36%):

- E-mail delivery isn't exactly guaranteed
- Should you buy a DVD player?
- Price low, demand high for Nintendo

Readers of top promoted stories (29%):

- 757 Crashes At Sea
- Israel, Palestinians Agree To Direct Talks
- Fuhrman Pleads Innocent To Perjury

MSNBC Story clusters

Readers of commerce and technology stories (36%):

- E-mail delivery isn't exactly guaranteed
- Should you buy a DVD player?
- Price low, demand high for Nintendo

Readers of top promoted stories (29%):

- 757 Crashes At Sea
- Israel, Palestinians Agree To Direct Talks
- Fuhrman Pleads Innocent To Perjury

Sports Readers (19%):

- Umps refusing to work is the right thing
- Cowboys are reborn in win over eagles
- Did Orioles spend money wisely?

MSNBC Story clusters

Readers of commerce and technology stories (36%):

- E-mail delivery isn't exactly guaranteed
- Should you buy a DVD player?
- Price low, demand high for Nintendo

Sports Readers (19%):

- Umps refusing to work is the right thing
- Cowboys are reborn in win over eagles
- Did Orioles spend money wisely?

Readers of top promoted stories (29%):

- 757 Crashes At Sea
- Israel, Palestinians Agree To Direct Talks
- Fuhrman Pleads Innocent To Perjury

Readers of “Softer” News (12%):

- The truth about what things cost
- Fuhrman Pleads Innocent To Perjury
- Real Astrology

Top 5 shows by user class

Class 1

- Power rangers
- Animaniacs
- X-men
- Tazmania
- Spider man

Class 2

- Young and restless
- Bold and the beautiful
- As the world turns
- Price is right
- CBS eve news

Class 3

- Tonight show
- Conan O'Brien
- NBC nightly news
- Later with Kinnear
- Seinfeld

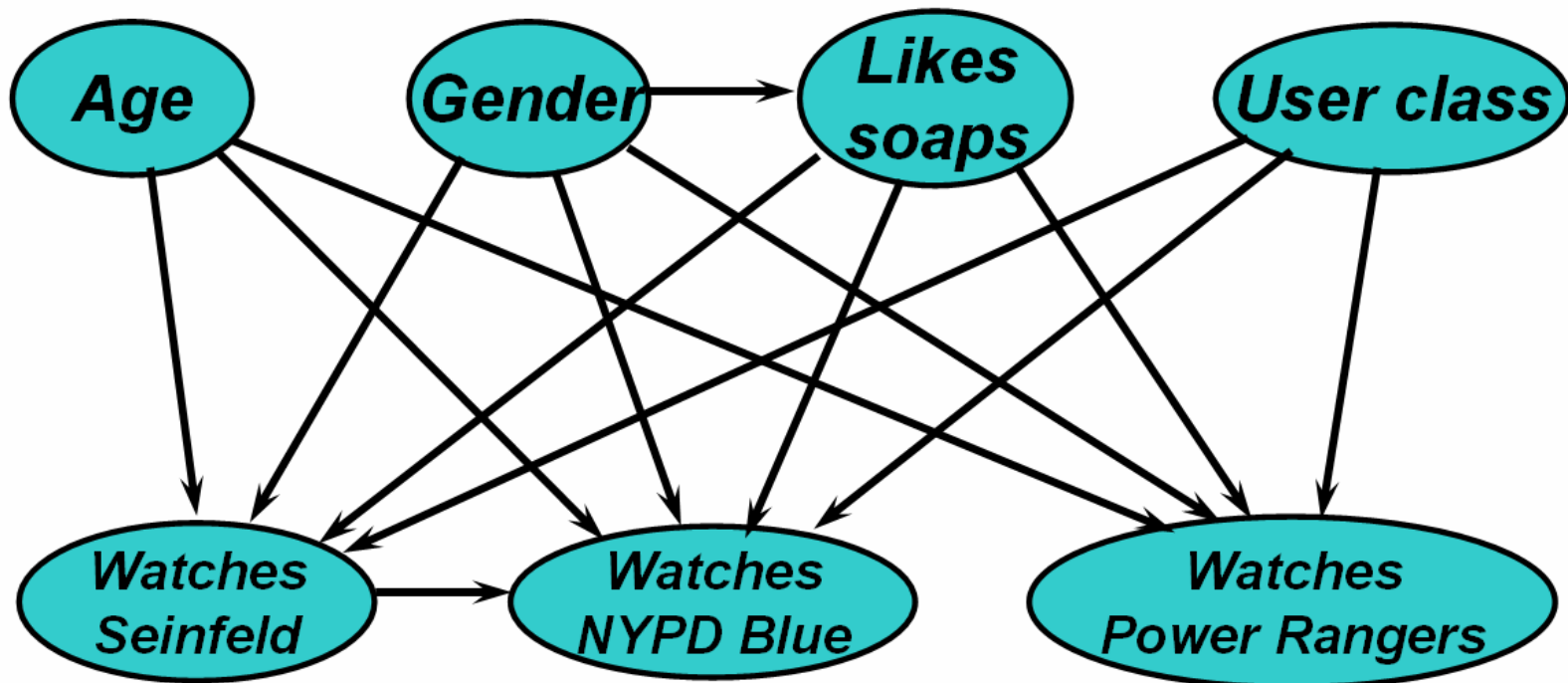
Class 4

- 60 minutes
- NBC nightly news
- CBS eve news
- Murder she wrote
- Matlock

Class 5

- Seinfeld
- Friends
- Mad about you
- ER
- Frasier

Richer model



What's old?

Decision theory & probability theory provide:

- principled models of belief and preference;
- techniques for:
 - ◆ integrating evidence (conditioning);
 - ◆ optimal decision making (max. expected utility);
 - ◆ targeted information gathering (value of info.);
 - ◆ parameter estimation from data.

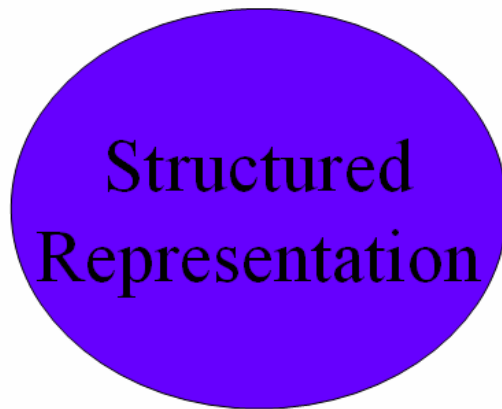
What's new?

What's new?

Bayesian networks exploit domain structure to allow compact representations of complex models.

What's new?

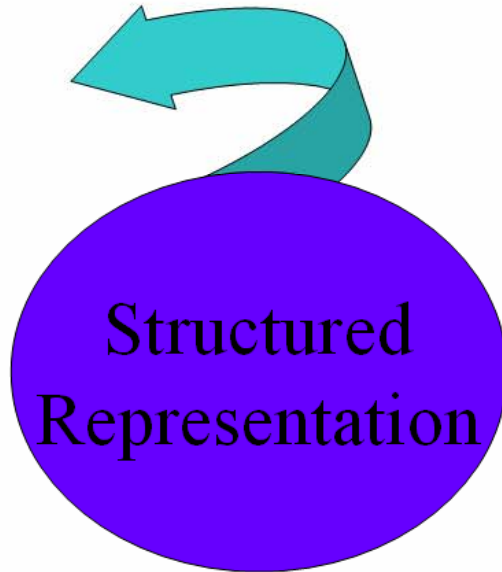
Bayesian networks exploit domain structure to allow compact representations of complex models.



What's new?

Bayesian networks exploit domain structure to allow compact representations of complex models.

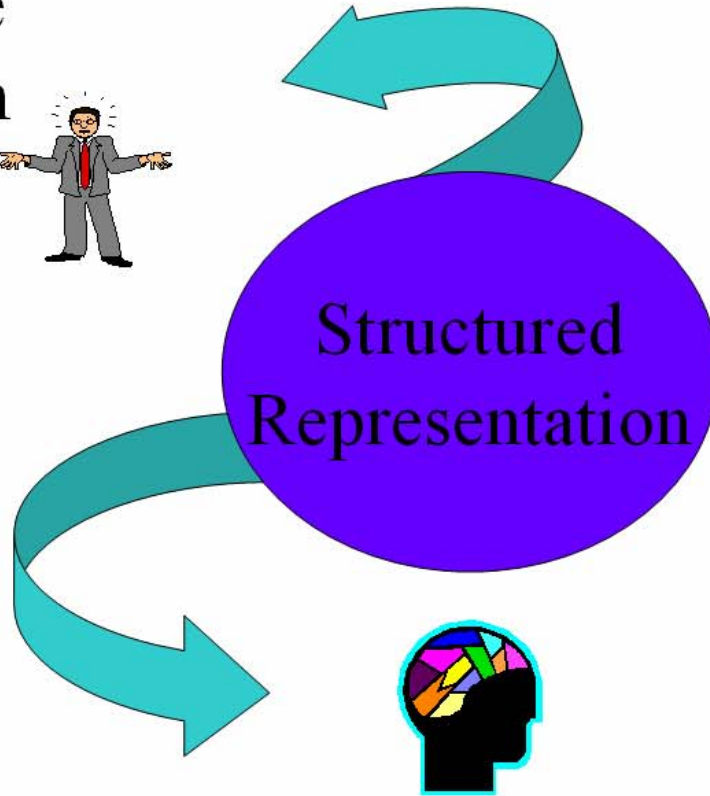
Knowledge
Acquisition



What's new?

Bayesian networks exploit domain structure to allow compact representations of complex models.

Knowledge
Acquisition



Inference

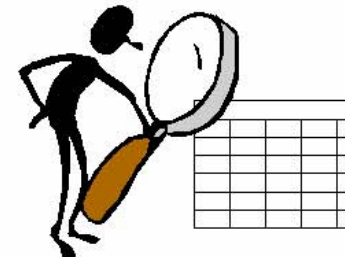
What's new?

Bayesian networks exploit domain structure to allow compact representations of complex models.

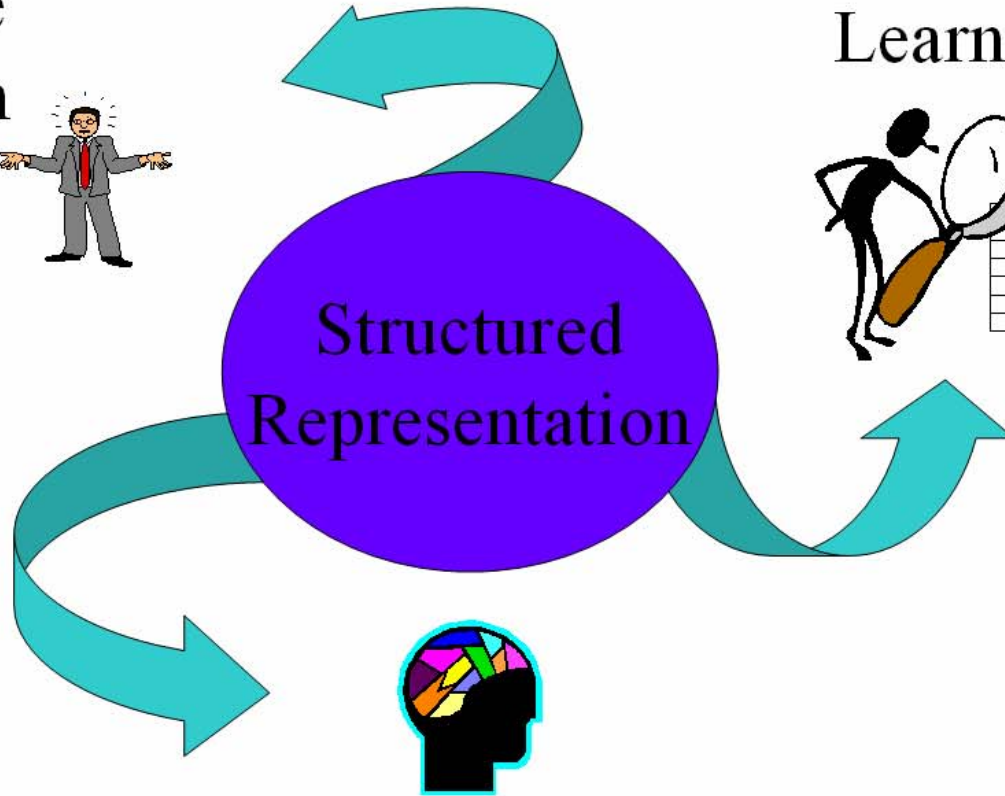
Knowledge
Acquisition



Learning



Structured
Representation



Inference

What's in our future?

What's in our future?

- Better models for:

What's in our future?

- Better models for:
 - ◆ preferences & utilities;

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.
- Inferring causality from data.

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.
- Inferring causality from data.
- More expressive representation languages:

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.
- Inferring causality from data.
- More expressive representation languages:
 - ◆ structured domains with multiple objects;

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.
- Inferring causality from data.
- More expressive representation languages:
 - ◆ structured domains with multiple objects;
 - ◆ levels of abstraction;

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.
- Inferring causality from data.
- More expressive representation languages:
 - ◆ structured domains with multiple objects;
 - ◆ levels of abstraction;
 - ◆ reasoning about time;

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.
- Inferring causality from data.
- More expressive representation languages:
 - ◆ structured domains with multiple objects;
 - ◆ levels of abstraction;
 - ◆ reasoning about time;
 - ◆ hybrid (continuous/discrete) models.

What's in our future?

- Better models for:
 - ◆ preferences & utilities;
 - ◆ not-so-precise numerical probabilities.
- Inferring causality from data.
- More expressive representation languages:
 - ◆ structured domains with multiple objects;
 - ◆ levels of abstraction;
 - ◆ reasoning about time;
 - ◆ hybrid (continuous/discrete) models.

