# Bayesian Networks and Decision-Theoretic Reasoning for Artificial Intelligence

Jack Breese

Microsoft Research

Daphne Koller

Stanford University

# Overview

■ Decision-theoretic techniques

  ◆ Explicit management of uncertainty and tradeoffs

  ◆ Probability theory

  ◆ Maximization of expected utility

■ Applications to AI problems

  ◆ Diagnosis

  ◆ Expert systems

  ◆ Planning

  ◆ Learning

# Science- AAAI-97

- Model Minimization in Markov Decision Processes

- Effective Bayesian Inference for Stochastic Programs

- Learning Bayesian Networks from Incomplete Data

- Summarizing CSP Hardness With Continuous Probability Distributions

- Speeding Safely: Multi-criteria Optimization in Probabilistic Planning

- Structured Solution Methods for Non-Markovian Decision Processes

# Applications

**@ COMPUTERWORLD**
*The online connection for information technology leaders*

**Microsoft's cost-cutting helps users**

**04/21/97**

**A Microsoft Corp. strategy to cut its support costs by letting users solve their own problems using electronic means is paying off for users.In March, the company began rolling out a series of Troubleshooting Wizards on its World Wide Web site.**
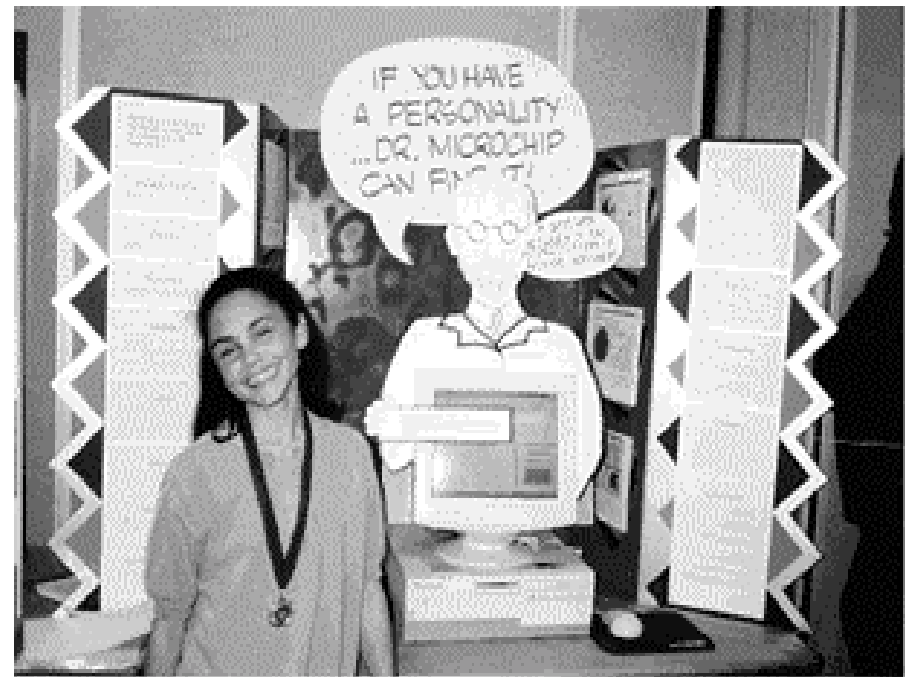
**Troubleshooting Wizards save time and money for users who don't have Windows NT specialists on hand at all times, said Paul Soares, vice president and general manager of Alden Buick Pontiac, a General Motors Corp. car dealership in Fairhaven, Mass**

# Teenage Bayes

**Microsoft Researchers Exchange Brainpower with Eighth-grader**

Teenager Designs Award-Winning Science Project

.. For her science project, which she called "Dr. Sigmund Microchip," Tovar wanted to create a computer program to diagnose the probability of certain personality types.  With only answers from a few questions, the program was able to accurately diagnose the correct personality type 90 percent of the time.



Elena Tovar stands proudly in front of "Dr. Sigmund Microchip," the science project she created using the advanced mathematical formulas that Microsoft Research uses to build artificial intelligence programs.

5

# Course Contents

» Concepts in Probability
  - ◆ Probability
  - ◆ Random variables
  - ◆ Basic properties (Bayes rule)
- Bayesian Networks
- Inference
- Decision making
- Learning networks from data
- Reasoning over time
- Applications

# Probabilities

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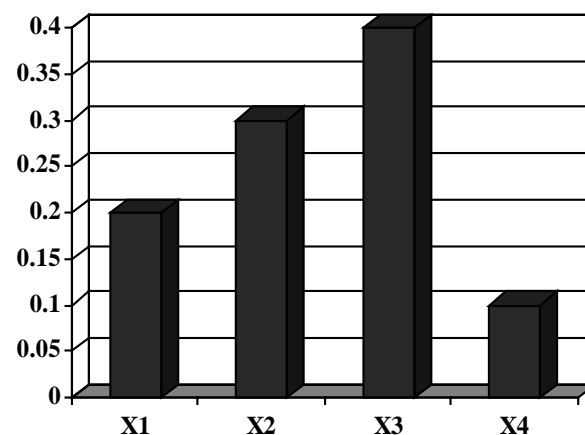    - $\xi$ is background state of information

# Discrete Random Variables

■ Finite set of possible outcomes

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

$$P(x_i) \geq 0$$



$$\sum_{i=1}^{n} P(x_i) = 1$$
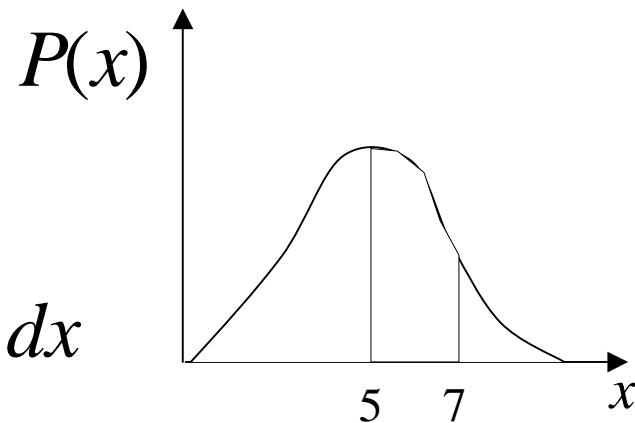
$X$ binary: $P(x) + P(\overline{x}) = 1$

# Continuous Random Variable

■ Probability distribution (density function) over continuous values

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

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

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

9

# More Probabilities

■ Conditional

$$P(x \mid y) \equiv P(X = x \mid Y = y)$$

◆ Probability that *X=x* given we know that *Y=y*

■ Joint

$$P(x, y) \equiv P(X = x \wedge Y = y)$$

◆ Probability that both *X=x* and Y=y

# Rules of Probability

- Product Rule

$$P(X,Y) = P(X \mid Y)P(Y) = P(Y \mid X)P(X)$$

- Marginalization

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

$X$ binary: $\quad P(Y) = P(Y, x) + P(Y, \overline{x})$

# Bayes Rule

$$P(H, E) = P(H \mid E)P(E) = P(E \mid H)P(H)$$

$$\boxed{P(H \mid E) = \frac{P(E \mid H)P(H)}{P(E)}}$$

$$P(h \mid e) = \frac{P(e \mid h)P(h)}{P(e,h) + P(e,\bar{h})}$$

$$= \frac{P(e \mid h)P(h)}{P(e \mid h)P(h) + P(e \mid \bar{h})P(\bar{h})}$$

# Course Contents

- ■ Concepts in Probability
  - » Bayesian Networks
    - ◆ Basics
    - ◆ Additional structure
    - ◆ Knowledge acquisition
- ■ Inference
- ■ Decision making
- ■ Learning networks from data
- ■ Reasoning over time
- ■ Applications

# Bayesian networks

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

14

# Bayesian Networks

$S \in \{no, light, heavy\}$ (Smoking) → (Cancer)

| 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⇓    C⇒* | *none* | *benign* | *malignant* |
|---|---|---|---|
| *no* | 0.768 | 0.024 | 0.008 |
| *light* | 0.132 | 0.012 | 0.006 |
| *heavy* | 0.035 | 0.010 | 0.005 |

# Marginalization

| $S\Downarrow$   $C\Rightarrow$ | none | benign | malig | total | |
|---|---|---|---|---|---|
| no | 0.768 | 0.024 | 0.008 | .80 | *P(Smoke)* |
| light | 0.132 | 0.012 | 0.006 | .15 | |
| heavy | 0.035 | 0.010 | 0.005 | .05 | |
| total | 0.935 | 0.046 | 0.019 | | |

*P(Cancer)*

# Bayes Rule Revisited

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

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

| Cancer= | none | benign | malignant |
|---|---|---|---|
| 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

19

# Independence

Age    Gender

*Age* and *Gender* are independent.

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

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

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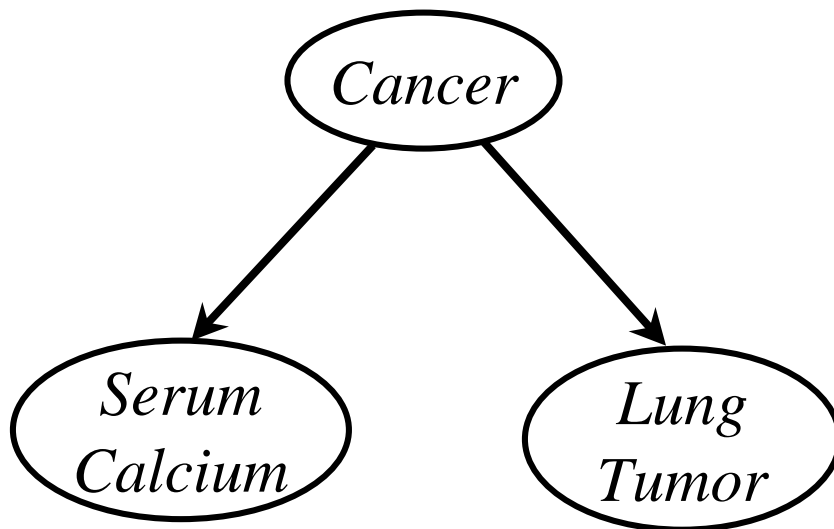$$P(A,G) = P(A/G)\, P(G) = P(A)P(G)$$

# Conditional Independence



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

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

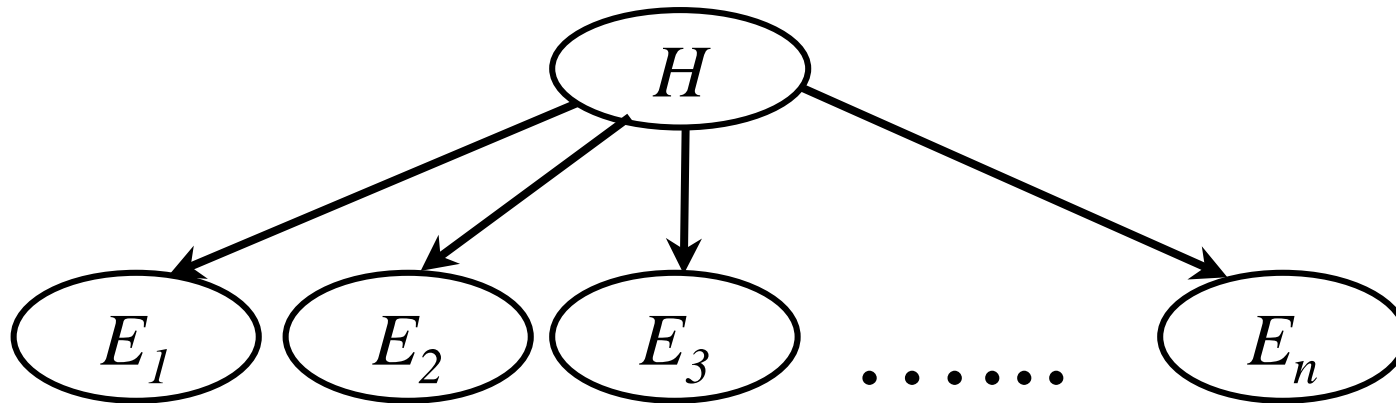# More Conditional Independence: Naïve Bayes

*Cancer*

*Serum Calcium*

*Lung Tumor*

*Serum Calcium* and *Lung Tumor* are dependent

*Serum Calcium* is independent of *Lung Tumor*, given *Cancer*
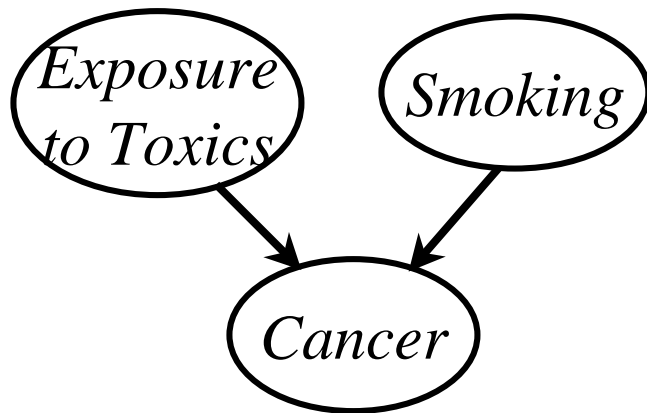
$$P(L/SC,C) = P(L/C)$$

# Naïve Bayes in general



$2n + 1$ *parameters:*
$$P(h)$$
$$P(e_i \mid h), P(e_i \mid \bar{h}), \ i = 1, \ldots, n$$

# More Conditional Independence: Explaining Away

Exposure
to Toxics
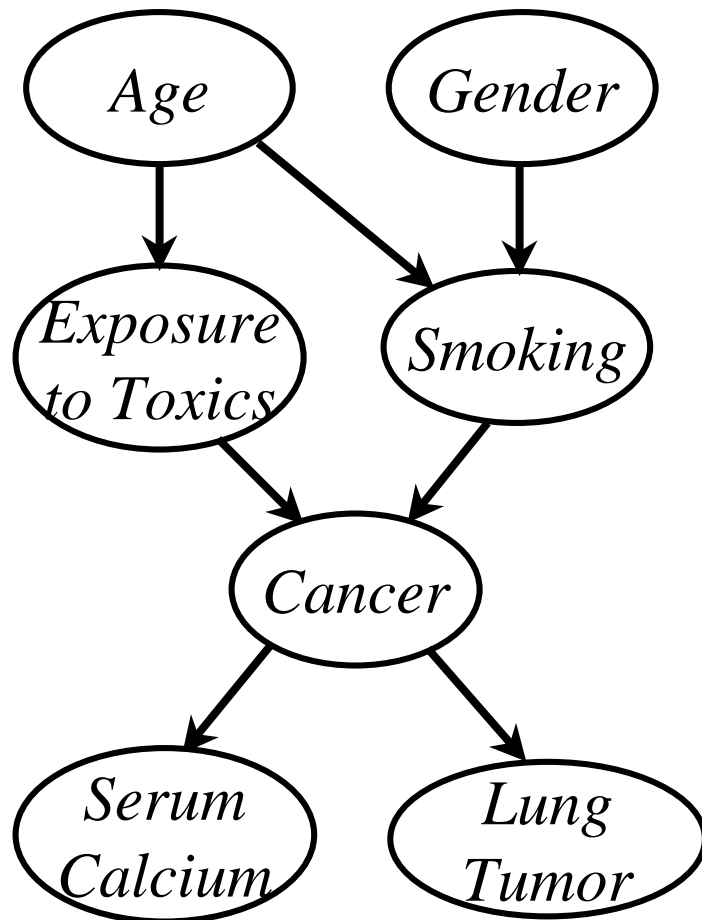
Smoking

Cancer

*Exposure to Toxics* and *Smoking* are independent

$$E \perp S$$

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

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

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

# Put it all together

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

$$P(A) \cdot P(G) \cdot$$



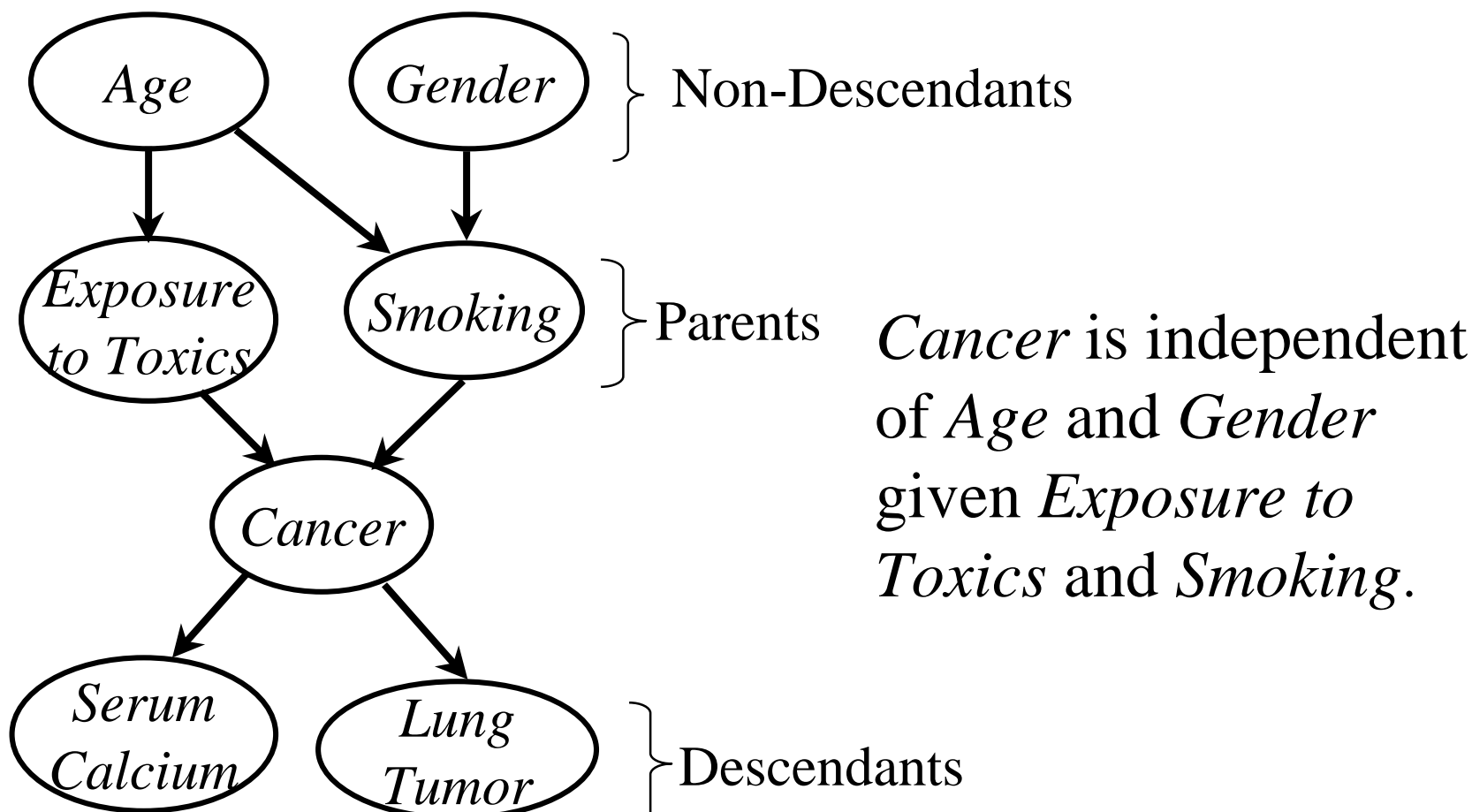$$P(E \mid A) \cdot P(S \mid A, G) \cdot$$

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

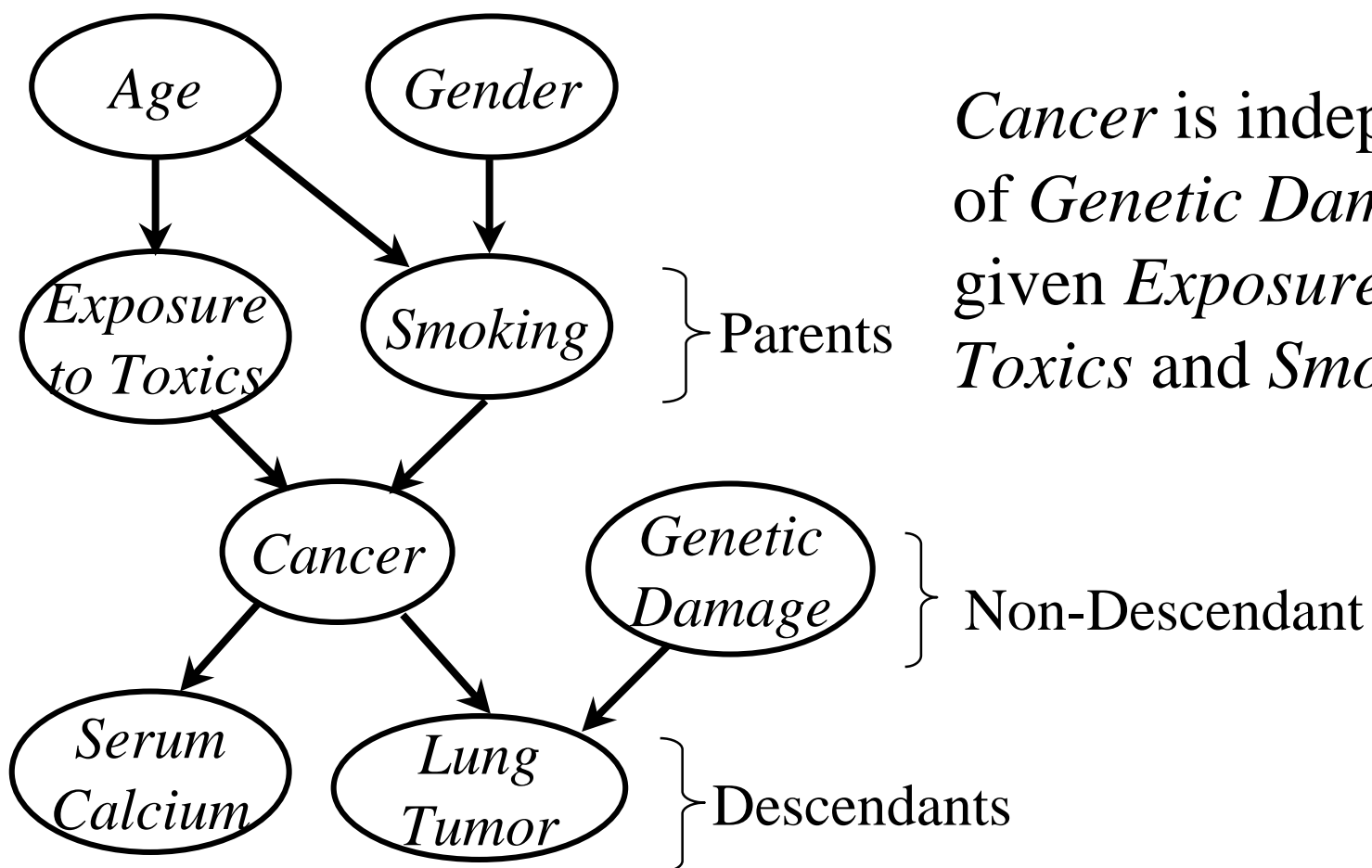$$P(SC \mid C) \cdot P(L \mid C)$$

# General Product (Chain) Rule for Bayesian Networks

$$P(X_1, X_2, \ldots, X_n) = \prod_{i=1}^{n} P(X_i \mid Pa_i)$$

$$Pa_i = parents(X_i)$$

# Conditional Independence

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

Age    Gender    } Non-Descendants

Exposure to Toxics    Smoking    } Parents

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

Cancer

Serum Calcium    Lung Tumor    } Descendants

# Another non-descendant

Age

Gender

Exposure to Toxics

Smoking

} Parents

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

Cancer

Genetic Damage

} Non-Descendant

Serum Calcium

Lung Tumor
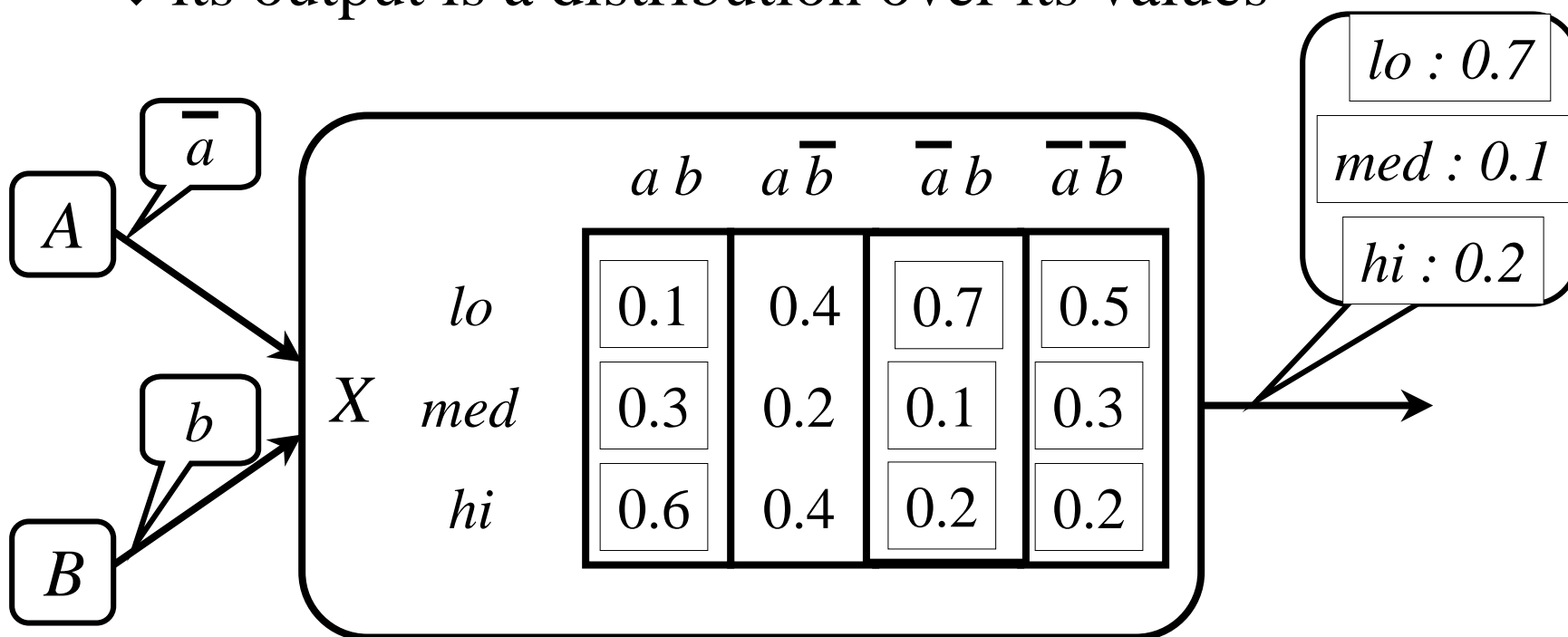
} Descendants

# 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.
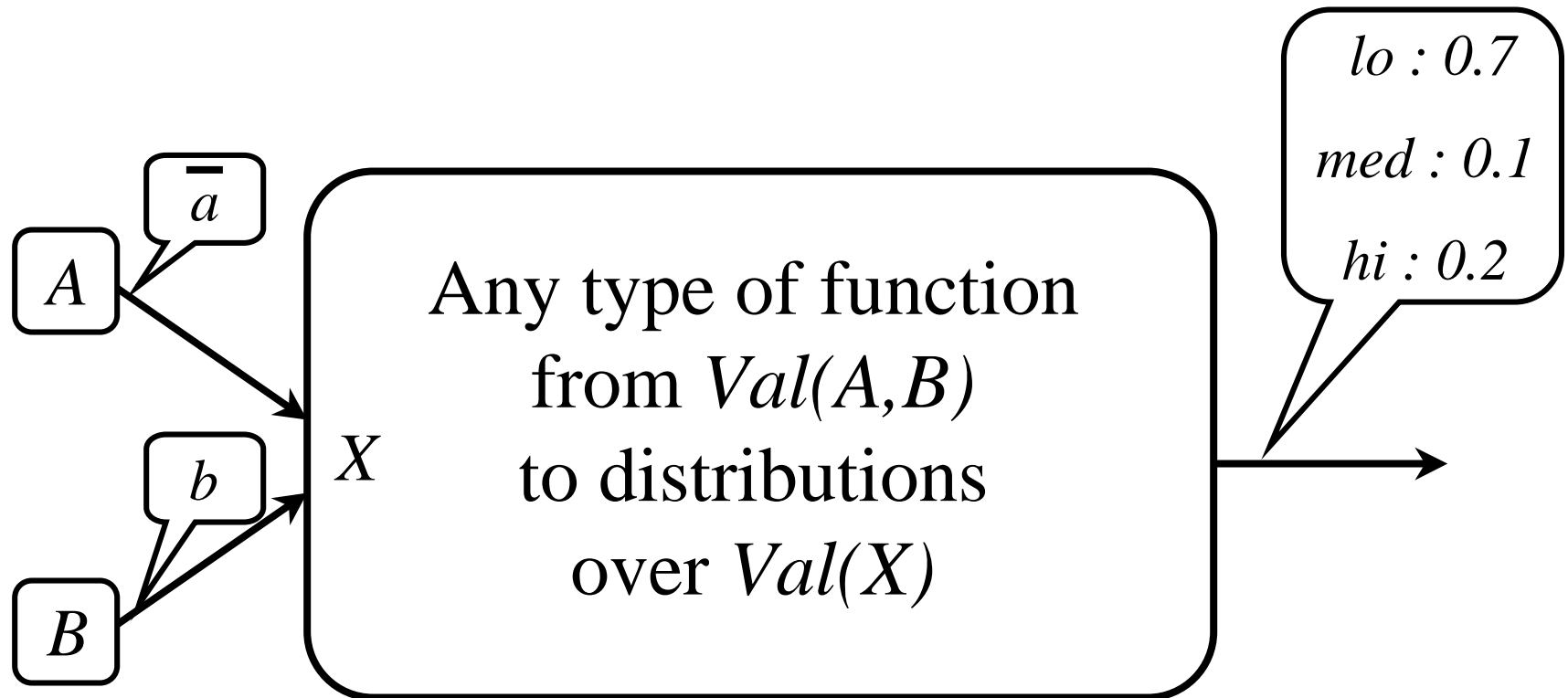
# Bayesian networks

■ Additional structure

- ◆ Nodes as functions

- ◆ Causal independence

- ◆ Context specific dependencies

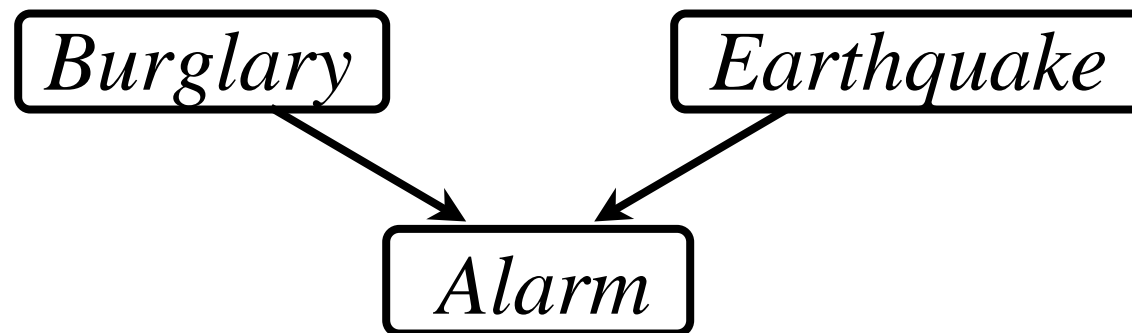- ◆ Continuous variables

- ◆ Hierarchy and model construction

# Nodes as functions

■ A BN node is conditional distribution function
   ◆ its parent values are the inputs
   ◆ its output is a distribution over its values

$\bar{a}$

A

$b$

B

X

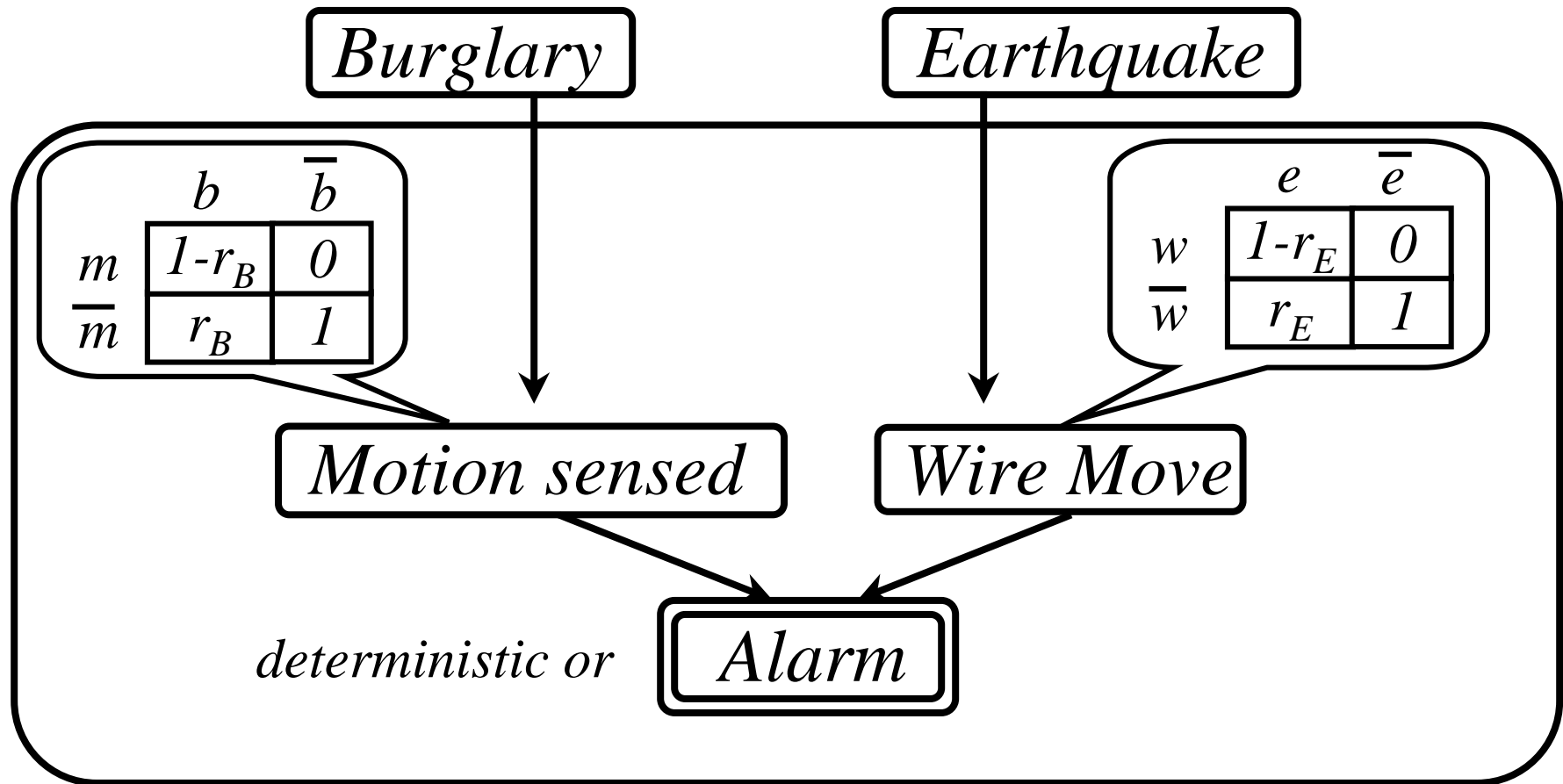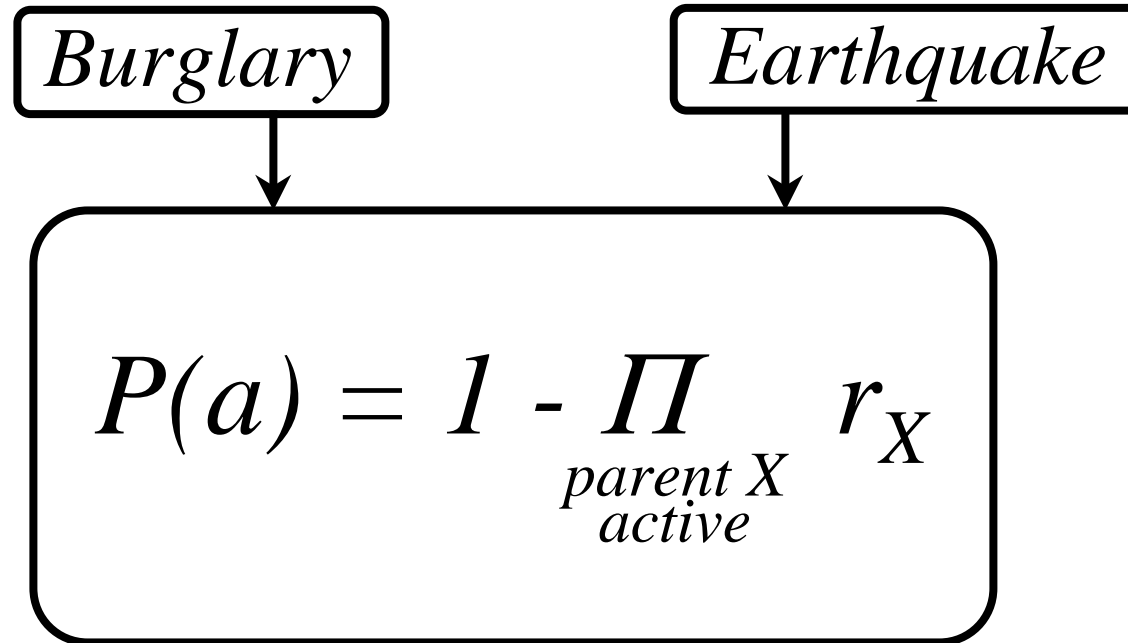| | $a\,b$ | $a\,\bar{b}$ | $\bar{a}\,b$ | $\bar{a}\,\bar{b}$ |
|---|---|---|---|---|
| lo | 0.1 | 0.4 | 0.7 | 0.5 |
| med | 0.3 | 0.2 | 0.1 | 0.3 |
| hi | 0.6 | 0.4 | 0.2 | 0.2 |

lo : 0.7

med : 0.1

hi : 0.2

# Causal Independence



■ *Burglary* causes *Alarm* iff motion sensor clear

■ *Earthquake* causes *Alarm* iff wire loose

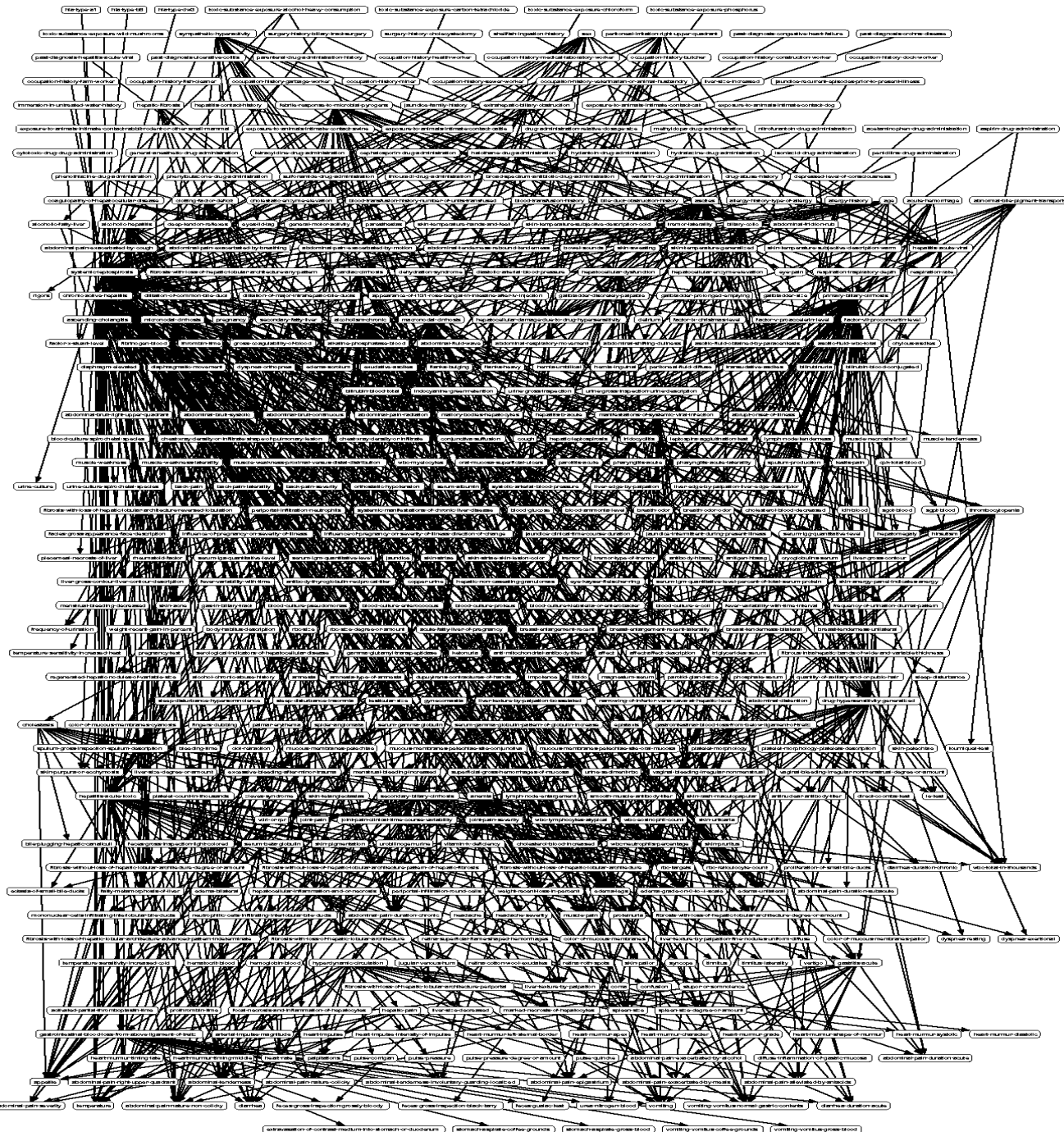■ Enabling factors are independent of each other

33

# Fine-grained model



| | $b$ | $\bar{b}$ |
|---|---|---|
| $m$ | $1-r_B$ | $0$ |
| $\bar{m}$ | $r_B$ | $1$ |

| | $e$ | $\bar{e}$ |
|---|---|---|
| $w$ | $1-r_E$ | $0$ |
| $\bar{w}$ | $r_E$ | $1$ |

**Burglary**

**Earthquake**

**Motion sensed**

**Wire Move**

*deterministic or*   **Alarm**

# Noisy-Or model

Alarm false only if all mechanisms independently inhibited

$$Burglary$$ $$Earthquake$$

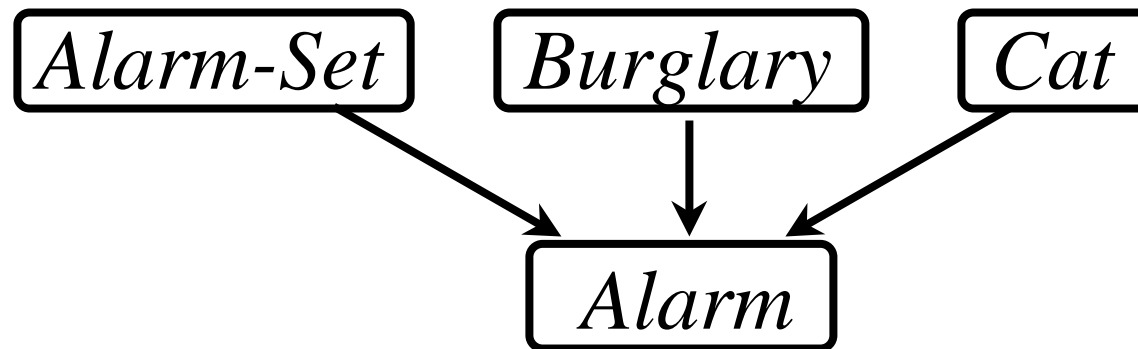$$P(a) = 1 - \prod_{\substack{parent\ X \\ active}} r_X$$

# of parameters is linear in the # of parents

35

# CPCS
# Network

36

# Context-specific Dependencies

```
┌─────────────┐   ┌─────────────┐   ┌───────┐
│  Alarm-Set  │   │  Burglary   │   │  Cat  │
└─────────────┘   └─────────────┘   └───────┘
          \             │             /
           \            │            /
            \           ▼           /
             ▶  ┌───────────────┐ ◀
                │     Alarm     │
                └───────────────┘
```
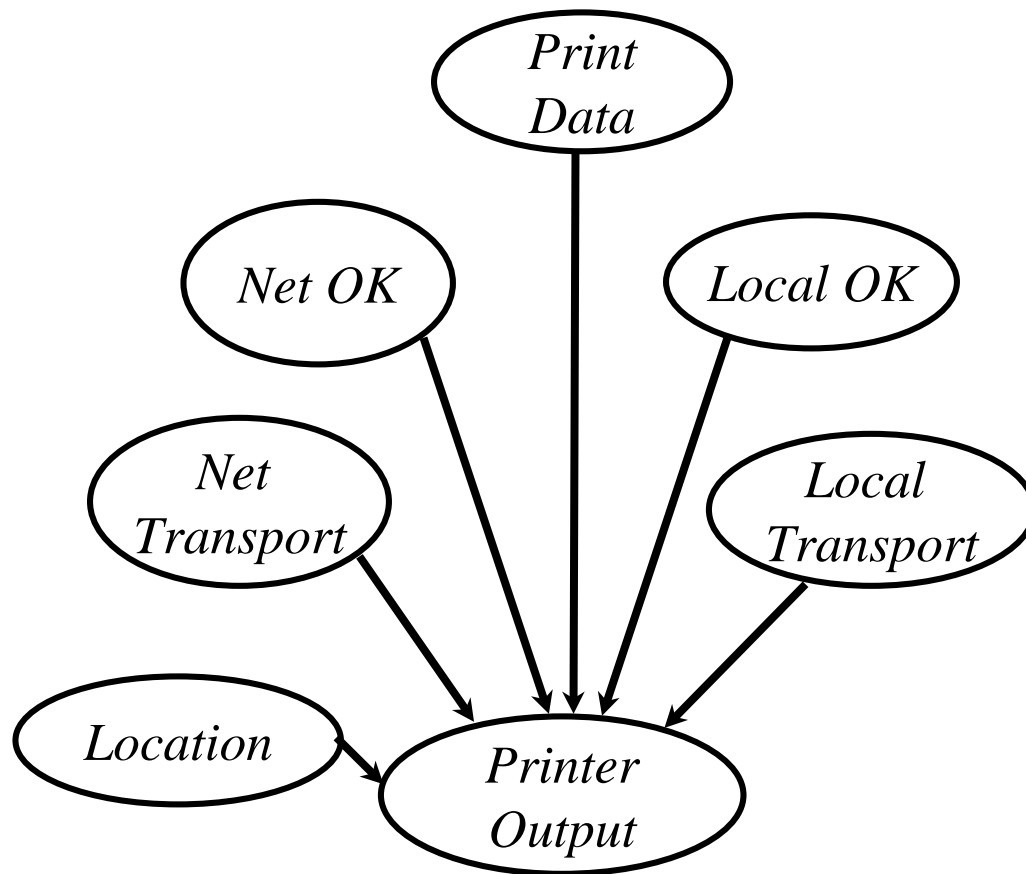
- *Alarm* can go off only if it is *Set*

- A burglar and the cat can both set off the alarm

- If a burglar comes in, the cat hides and does not set off the alarm
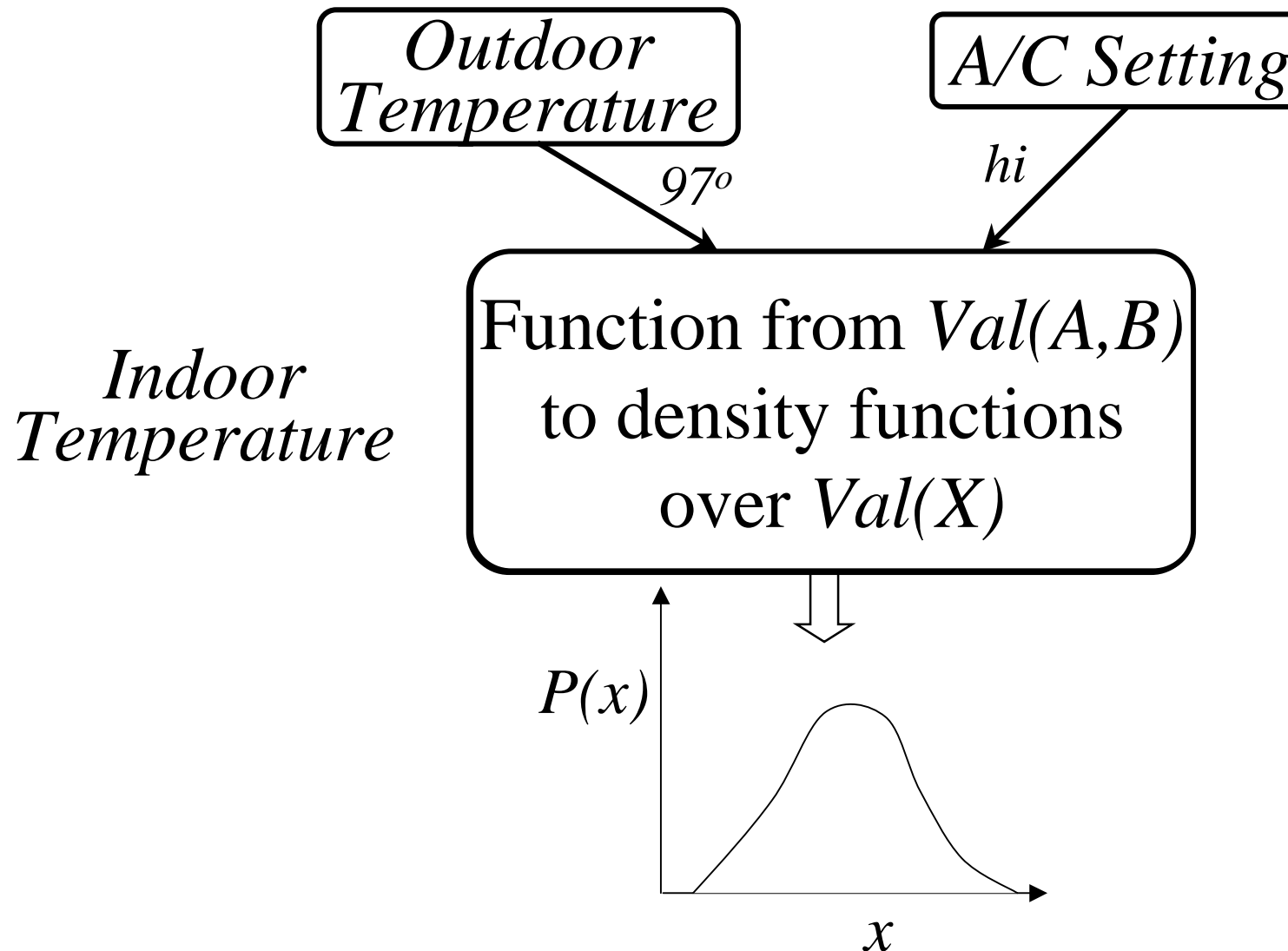
# Asymmetric dependencies

$\boxed{\textit{Alarm-Set}}$ $\boxed{\textit{Burglary}}$ $\boxed{\textit{Cat}}$

Node function
represented
as a tree

$A$ $\overline{s}$ $S$ $s$

$(a: 0, \overline{a} : 1)$ $\overline{b}$ $B$ $b$

$\overline{c}$ $C$ $c$ $(a: 0.9, \overline{a} : 0.1)$

$(a: 0.01, \overline{a} : 0.99)$ $(a: 0.6, \overline{a} : 0.4)$

- *Alarm* independent of
  - ◆ *Burglary*, *Cat* given $\overline{s}$
  - ◆ *Cat* given $s$ and $b$

38

# Asymmetric Assessment

Print Data

Net OK

Local OK

Net Transport

Local Transport

Location

Printer Output

**Asymmetric Assessment of Printer Output**

Assessment Hierarchy:

- PrintDataOut
  - Yes
    - Printer Location
      - Local
        - LOCAL Transport
          - Yes
            - Local Printer OK
              - Normal
              - Abnormal
          - No
      - Network
        - NET Transport
          - Yes
            - Net Printer OK
              - Normal
              - Abnormal
          - No
  - No

# Continuous variables

Outdoor Temperature

A/C Setting

97°

hi

Indoor Temperature

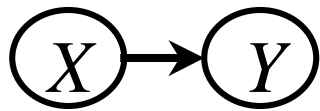Function from *Val(A,B)* to density functions over *Val(X)*

*P(x)*

*x*

# Gaussian (normal) distributions

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



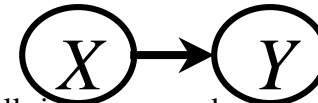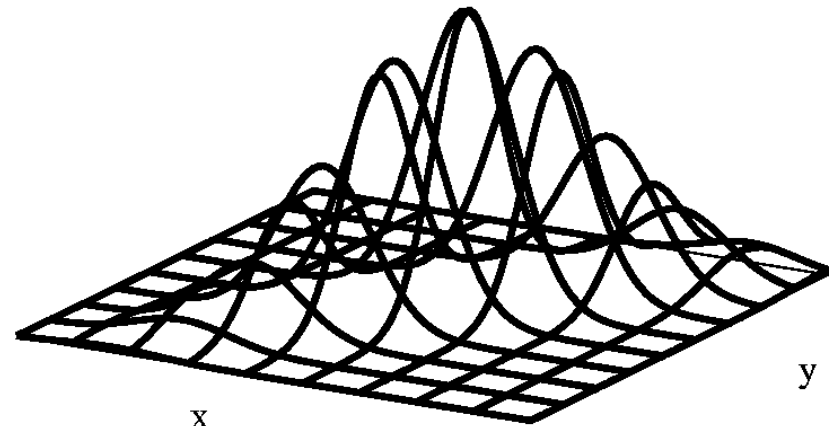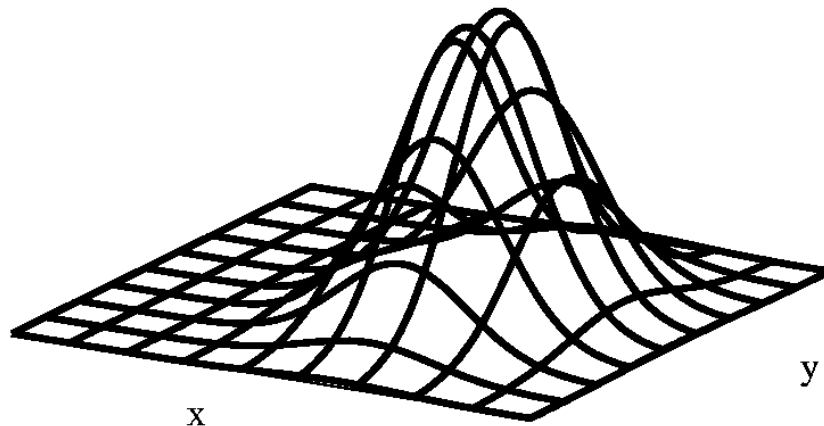*different mean*   *different variance*

41

# Gaussian networks

$$X \sim N(\mu, \sigma_X^2)$$

$X \rightarrow Y$

$$Y \sim N(ax + b, \sigma_Y^2)$$

Each variable is a linear function of its parents, with Gaussian noise

Joint probability density functions:



y

x

$X$   $Y$



y

x

$X \rightarrow Y$

# Composing functions

■ Recall: a BN node is a function

■ We can compose functions to get more complex functions.

■ The result: A hierarchically structured BN.

■ Since functions can be called more than once, we can reuse a BN model fragment in multiple contexts.
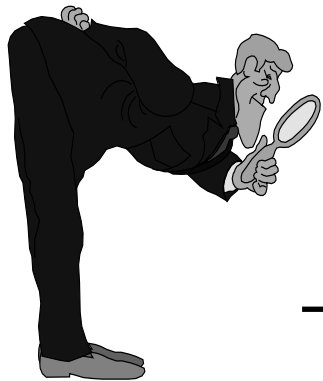
43

# Bayesian Networks

■ Knowledge acquisition

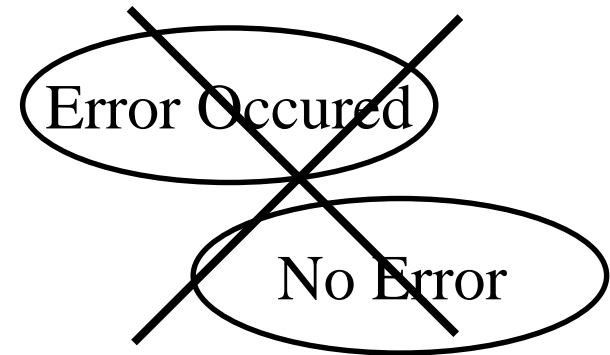    ◆ Variables

    ◆ Structure

    ◆ Numbers

# What is a variable?
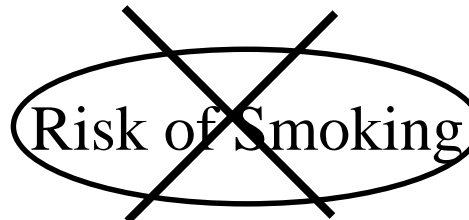
■ Collectively exhaustive, mutually exclusive values

$$x_1 \vee x_2 \vee x_3 \vee x_4$$

$$\neg(x_i \wedge x_j) \quad i \neq j$$

Error Occured

No Error

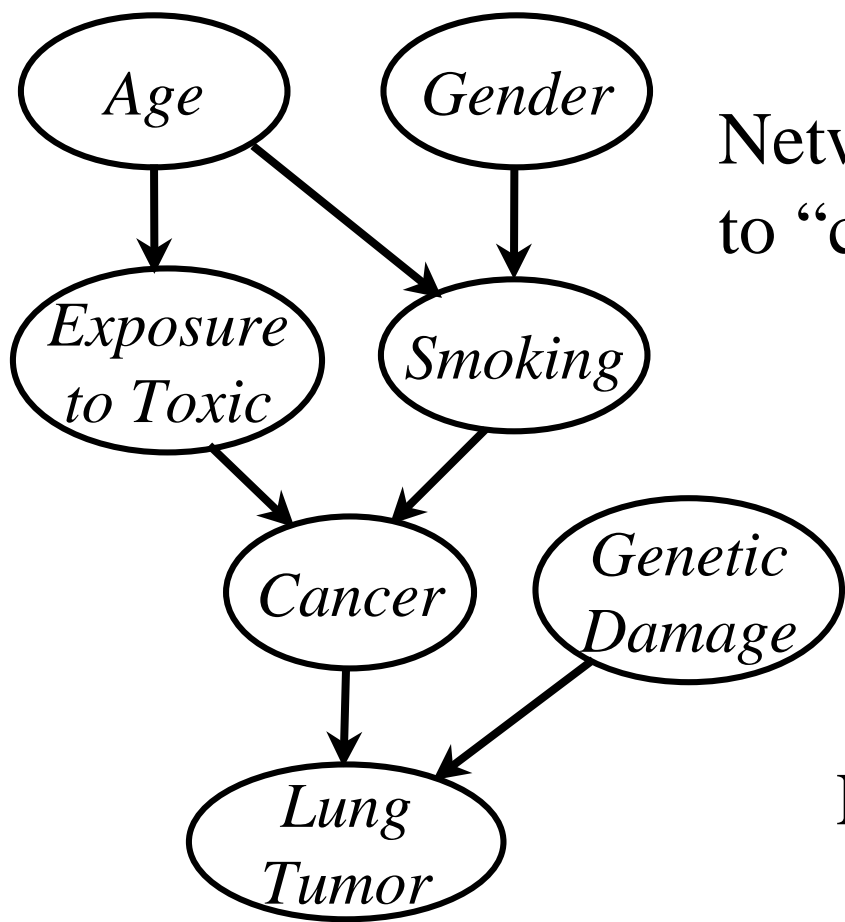■ Values versus Probabilities

Risk of Smoking

Smoking

46

# Clarity Test:
# Knowable in Principle

■ Weather  {Sunny, Cloudy, Rain, Snow}

■ Gasoline: Cents per gallon

■ Temperature { ≥ 100F , < 100F}

■ User needs help on Excel Charting {Yes, No}

■ User's personality {dominant, submissive}

# Structuring

Age

Gender

Exposure to Toxic

Smoking

Cancer

Genetic Damage

Lung Tumor

Network structure corresponding to "causality" is usually good.

Extending the conversation.

# Do the numbers really matter?

■ Second decimal usually does not matter

■ Relative Probabilities

**Assess probabilities for: I-TypingSpeed_avg**

**I-TypingSpeed**

| E-Arousal | Fast | Normal | Slow |
|-----------|------|--------|------|
| Passive | .20 | .28 | .52 |
| Neutral | .33 | .33 | .33 |
| Excited | .56 | .27 | .16 |

Ok   Cancel

■ Zeros and Ones

■ Order of Magnitude : $10^{-9}$ vs $10^{-6}$

■ Sensitivity Analysis

# Bayesian Networks and Structure


Windows '95 Print Troubleshooter -- MINIPTS.DSC

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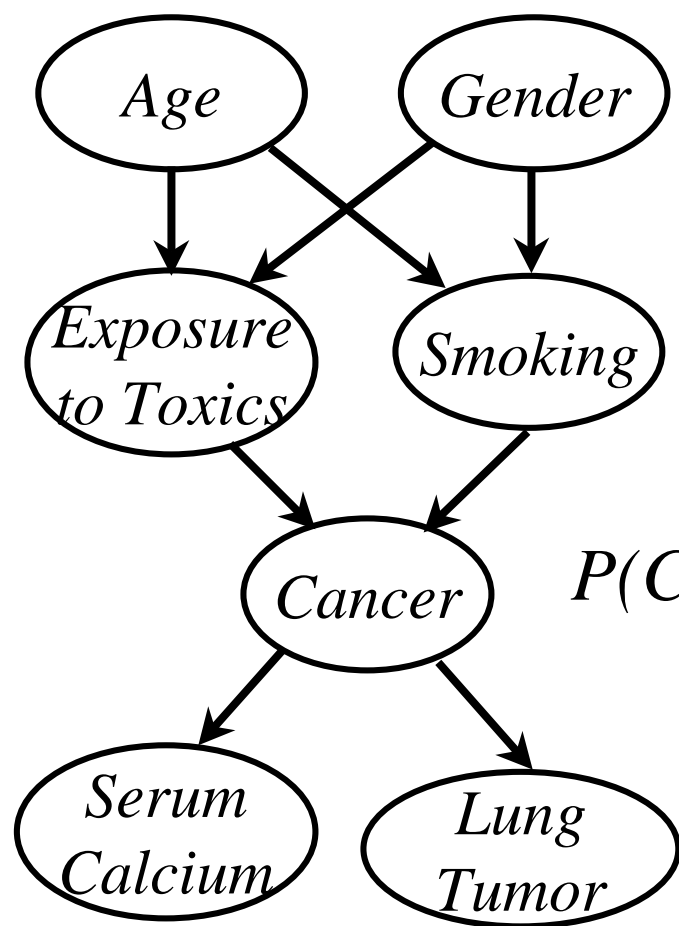  - 145 to 55 for a small hardware network;
  - 133,931,430 to 8254 for CPCS !!

# 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
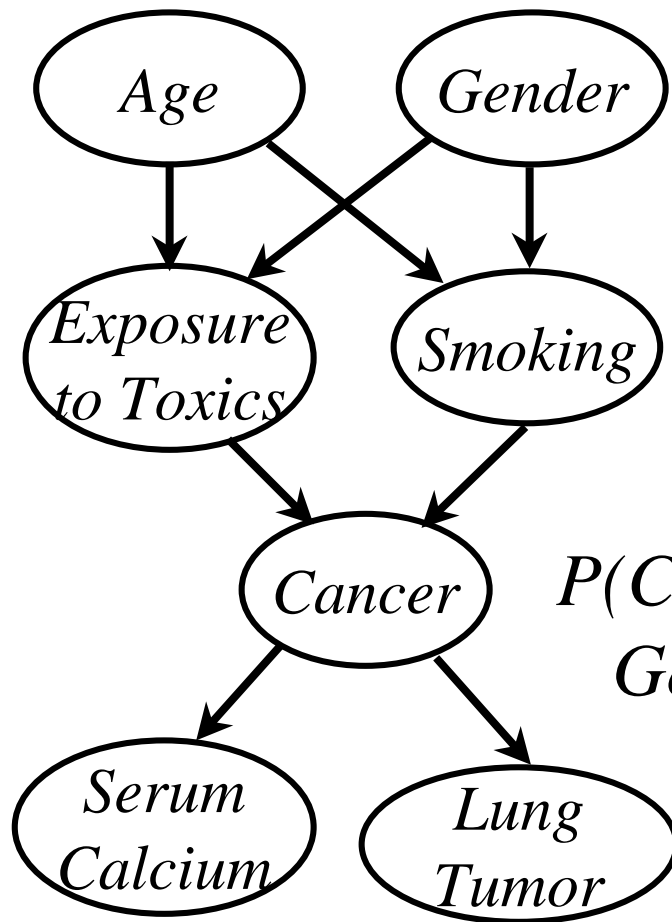
- ■ Exploiting structure

- ■ Approximate inference

# Predictive Inference



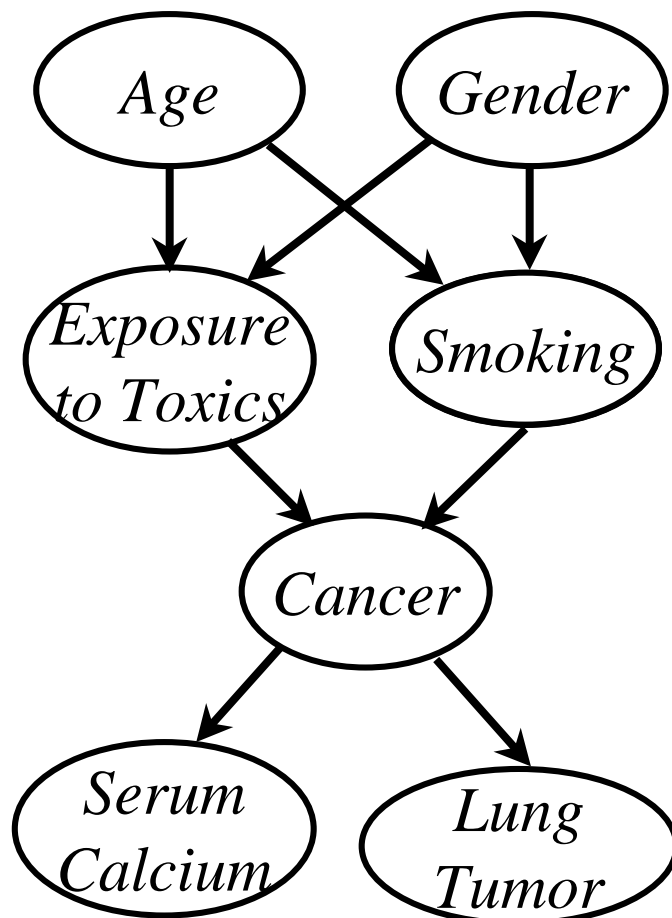How likely are elderly males to get malignant cancers?

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

53

# Combined



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

*P(C=malignant | Age>60, Gender= male, Serum Calcium = high)*

# Explaining away

Age → Exposure to Toxics

Gender → Smoking

Age → Smoking

Gender → Exposure to Toxics

Exposure to Toxics → Cancer

Smoking → Cancer

Cancer → Serum Calcium

Cancer → Lung Tumor

- ■ 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
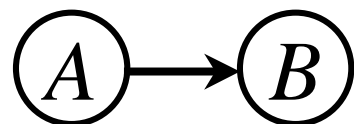
■ Find *P(Q=q/E= e)*

　◆ *Q*  the query variable

　◆ *E*  set of evidence variables

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

$X_1, \ldots, X_n$ are network variables except $Q, \boldsymbol{E}$

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

# Basic Inference



$P(b) = ?$

# Product Rule

$S \rightarrow C$

■ *P(C,S) = P(C|S) P(S)*

| $S\Downarrow$    $C\Rightarrow$ | *none* | *benign* | *malignant* |
|---|---|---|---|
| *no* | 0.768 | 0.024 | 0.008 |
| *light* | 0.132 | 0.012 | 0.006 |
| *heavy* | 0.035 | 0.010 | 0.005 |

# Marginalization

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

*P(Smoke)*

*P(Cancer)*

# Basic Inference

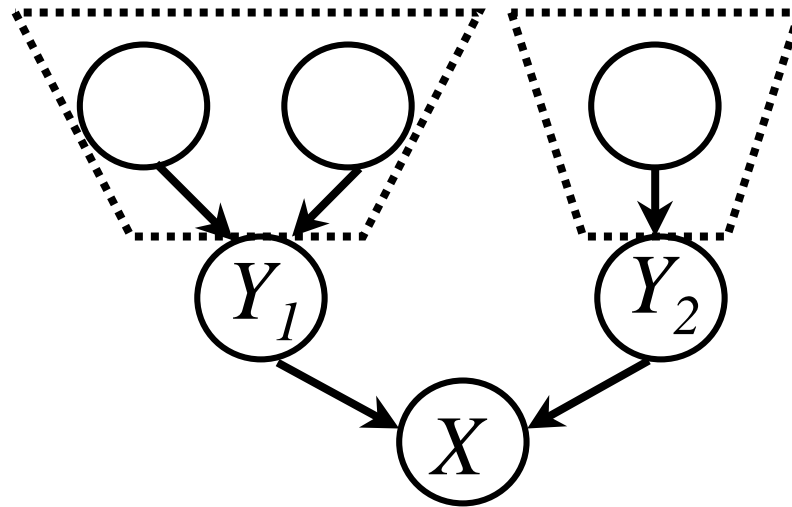$$A \rightarrow B \rightarrow C$$

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

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

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

$$= \sum_b P(c \mid b) \underbrace{\sum_a P(b \mid a) \, P(a)}_{P(b)}$$
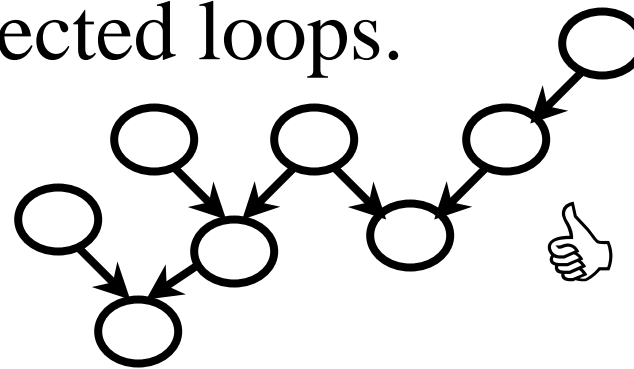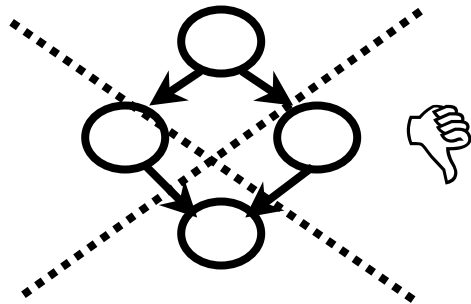
# Inference in trees



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

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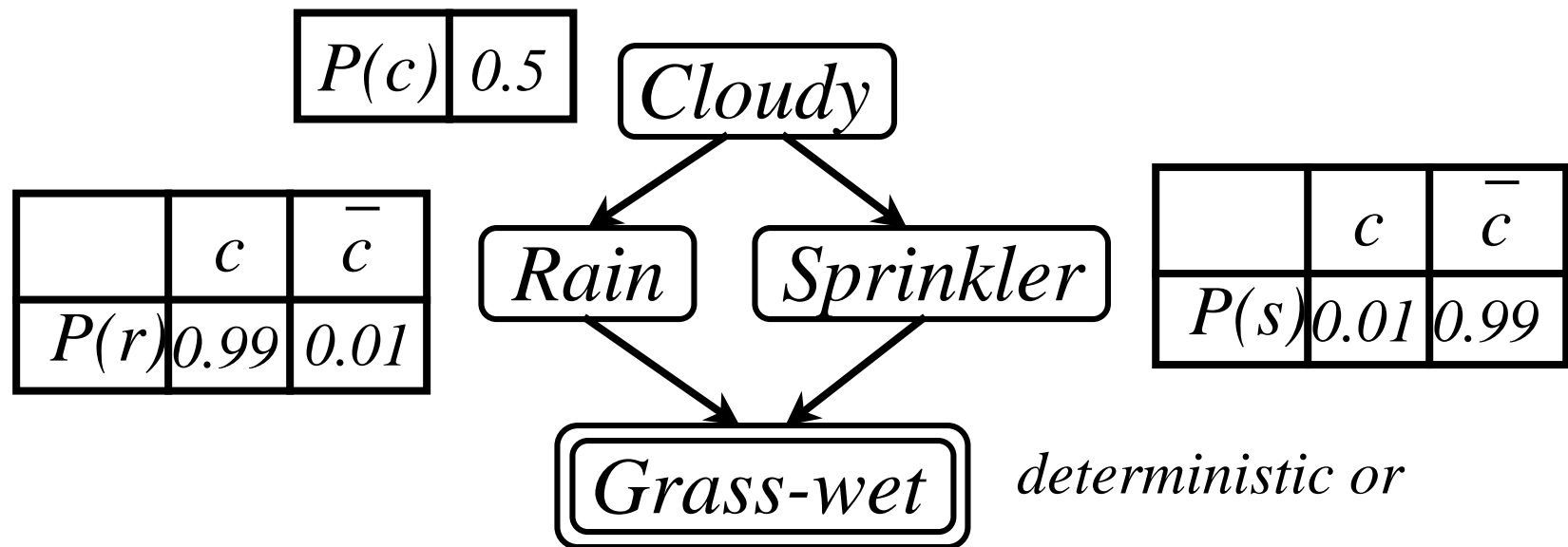
**Theorem:** Inference in a singly connected network can be done in linear time*.

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

* in network size including table sizes.

# The problem with loops

| P(c) | 0.5 |
|------|-----|

**Cloudy**

**Rain**    **Sprinkler**

|  | $c$ | $\bar{c}$ |
|------|------|------|
| P(r) | 0.99 | 0.01 |

|  | $c$ | $\bar{c}$ |
|------|------|------|
| P(s) | 0.01 | 0.99 |

**Grass-wet**    *deterministic or*

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.

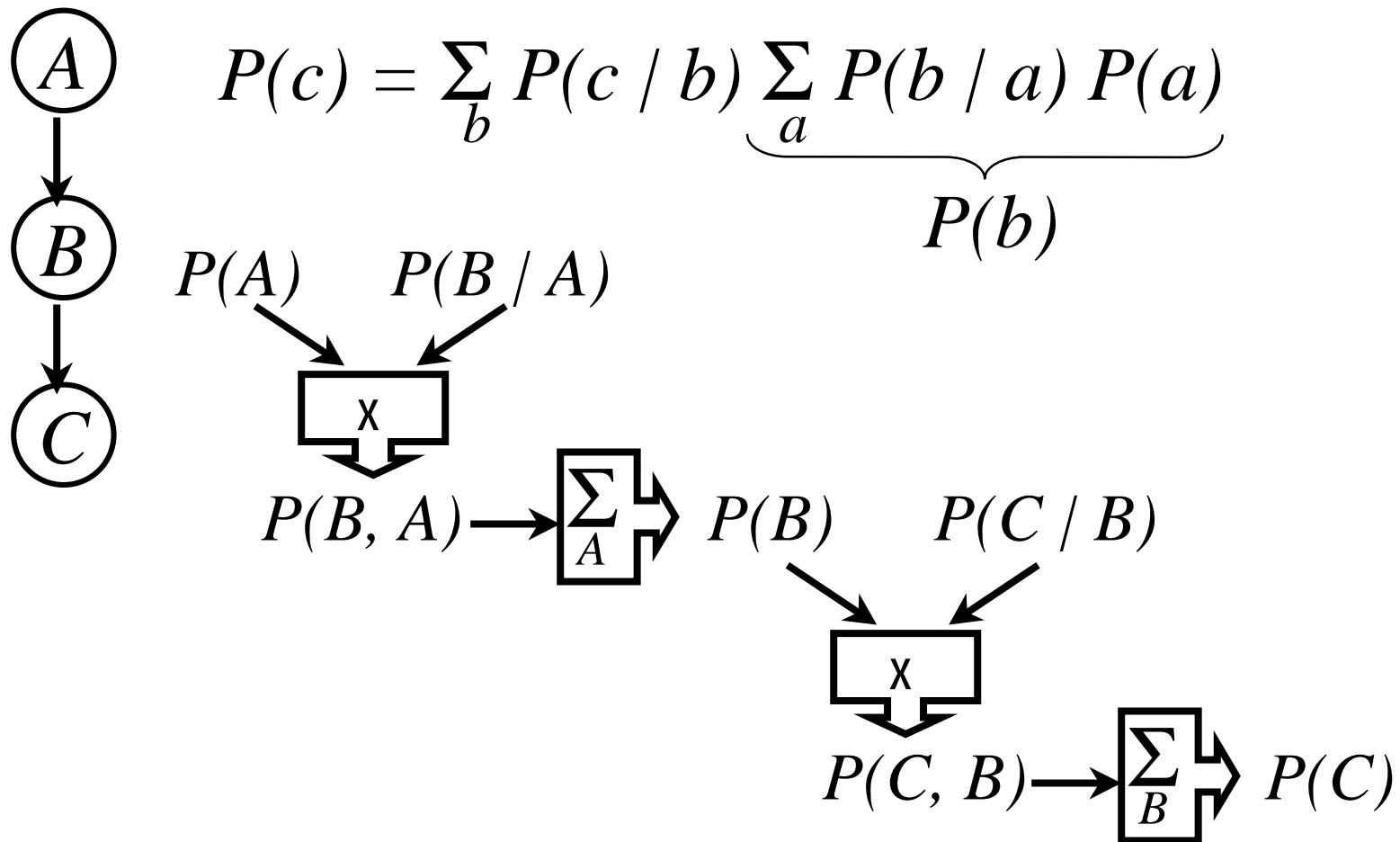$$P(\bar{g}) = \overbrace{P(\bar{g}\,/\,r,\,s)}^{0} P(r,\,s) + \overbrace{P(\bar{g}\,/\,r,\,\bar{s})}^{0} P(r,\,\bar{s})$$

$$+ \underbrace{P(\bar{g}\,/\,\bar{r},\,s)}_{0} P(\bar{r},\,s) + \underbrace{P(\bar{g}\,/\,\bar{r},\,\bar{s})}_{1} P(\bar{r},\,\bar{s})$$

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

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

problem

# Variable elimination



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

$P(A)$   $P(B \mid A)$

X

$P(B, A) \longrightarrow \underset{A}{\Sigma} \quad P(B) \qquad P(C \mid B)$
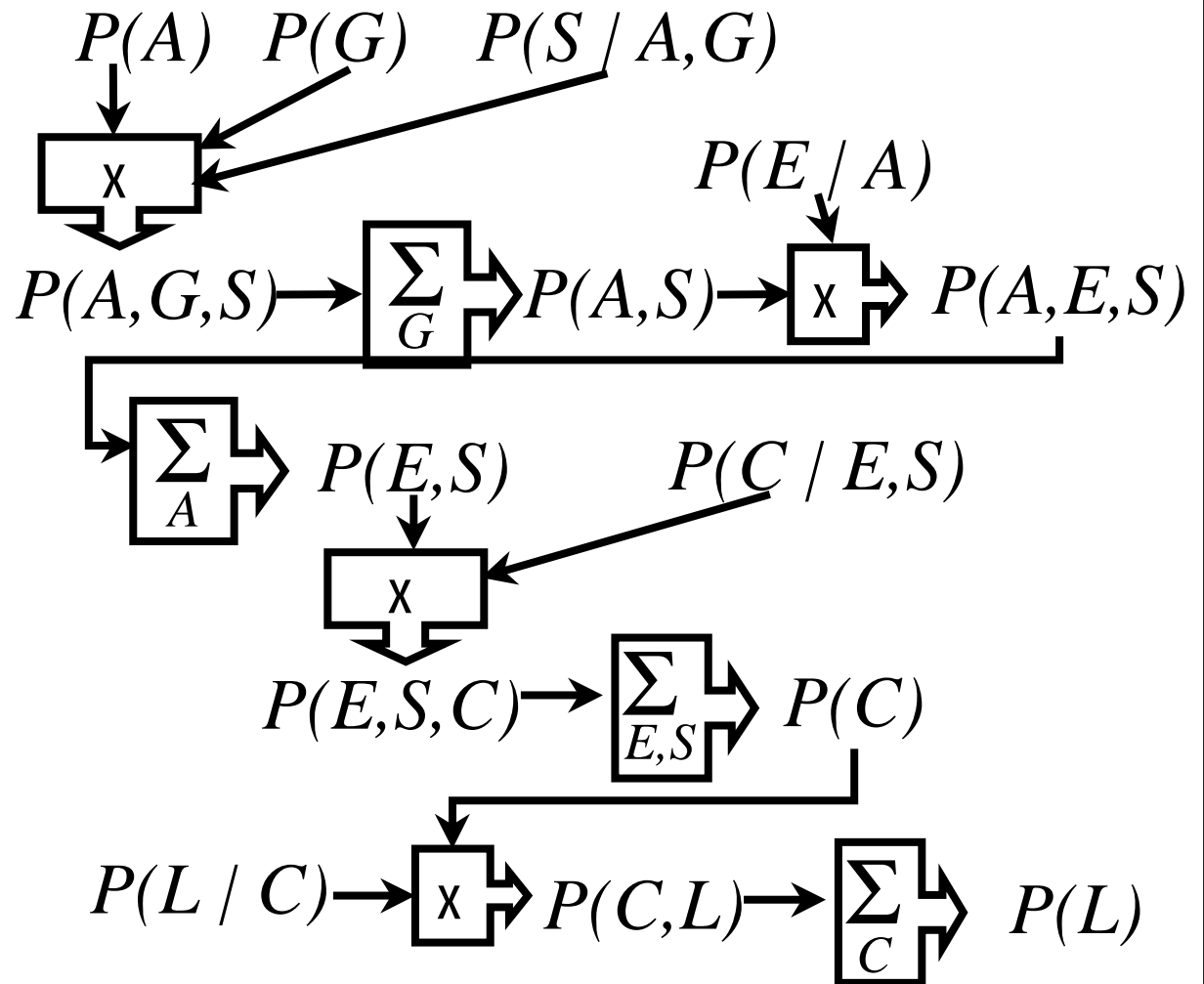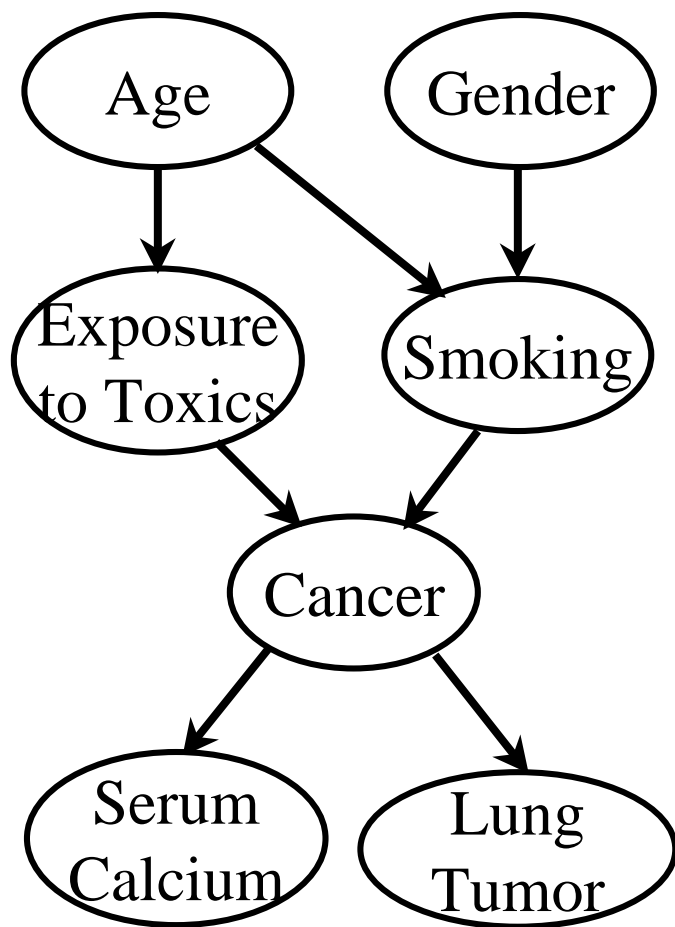
X

$P(C, B) \longrightarrow \underset{B}{\Sigma} \quad P(C)$

65

# 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.
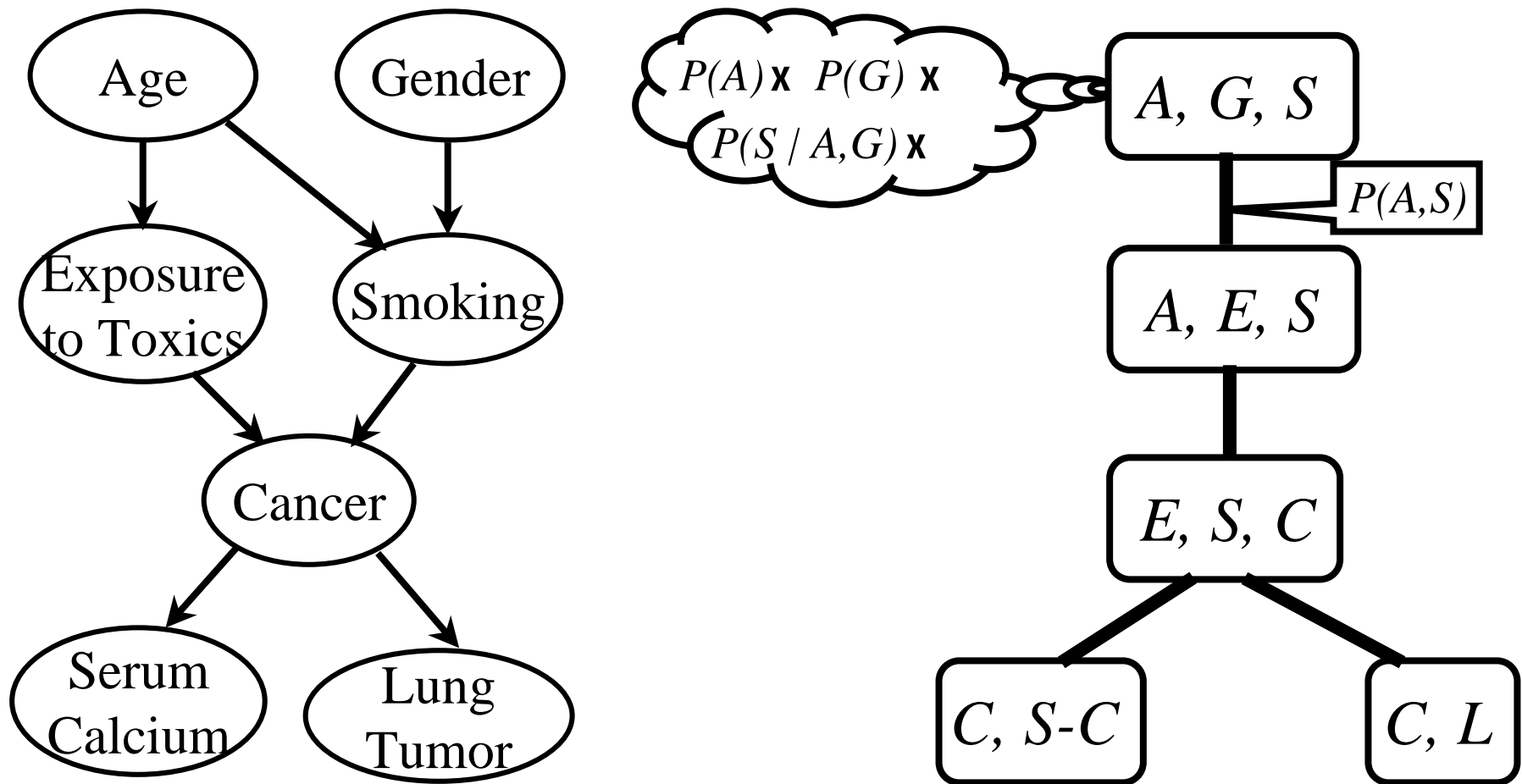
# Variable Elimination with loops



Complexity is exponential in the size of the factors

# Join trees*

A join tree is a partially precompiled factorization



Age

Gender

Exposure to Toxics

Smoking

Cancer

Serum Calcium

Lung Tumor

$P(A) \times P(G) \times P(S \mid A,G) \times$
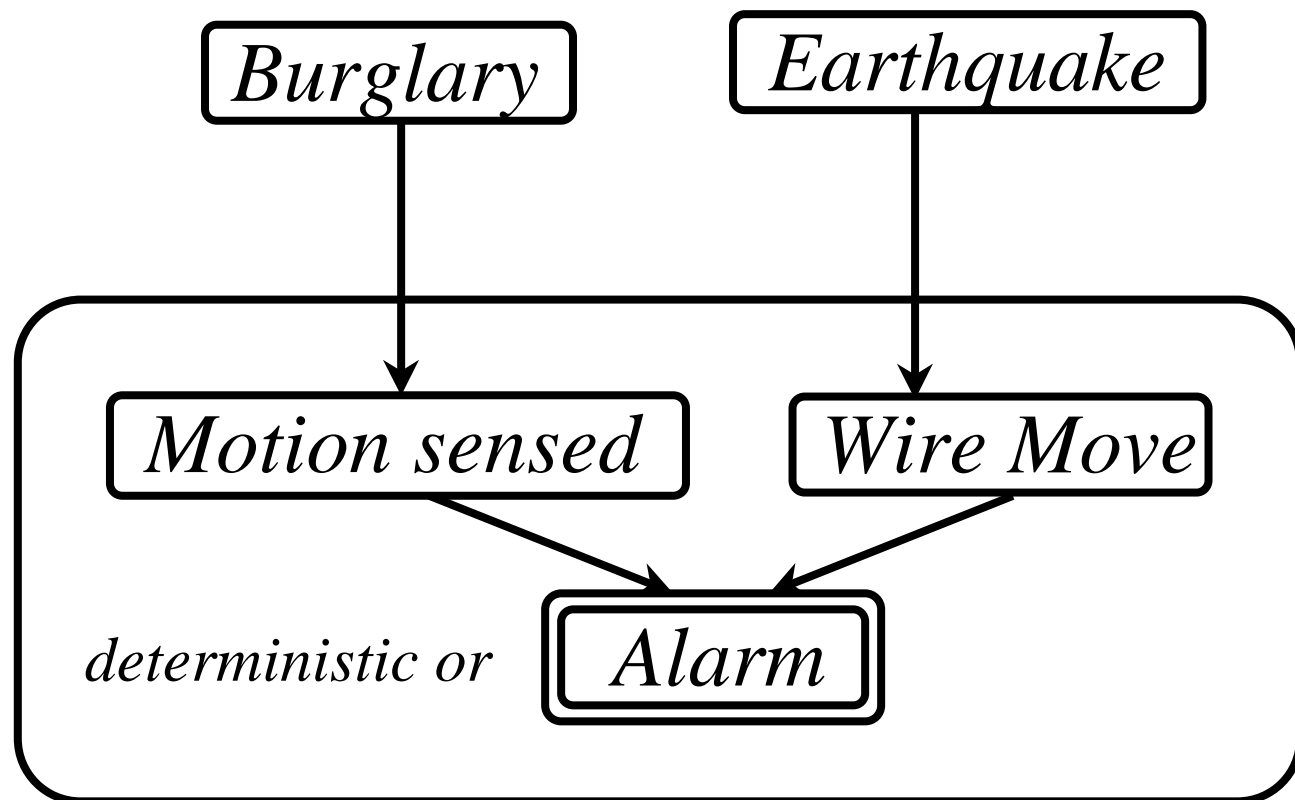
$A, G, S$

$P(A,S)$

$A, E, S$

$E, S, C$

$C, S\text{-}C$

$C, L$

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

# Exploiting Structure

Idea: explicitly decompose nodes

Noisy or:



Burglary

Earthquake

Motion sensed

Wire Move

*deterministic or*

Alarm

# Noisy-or decomposition



**Smaller families**
⇨ **Smaller factors**
⇨ **Faster inference**

# Inference with continuous variables

- Gaussian networks: polynomial time inference regardless of network structure

- Conditional Gaussians:
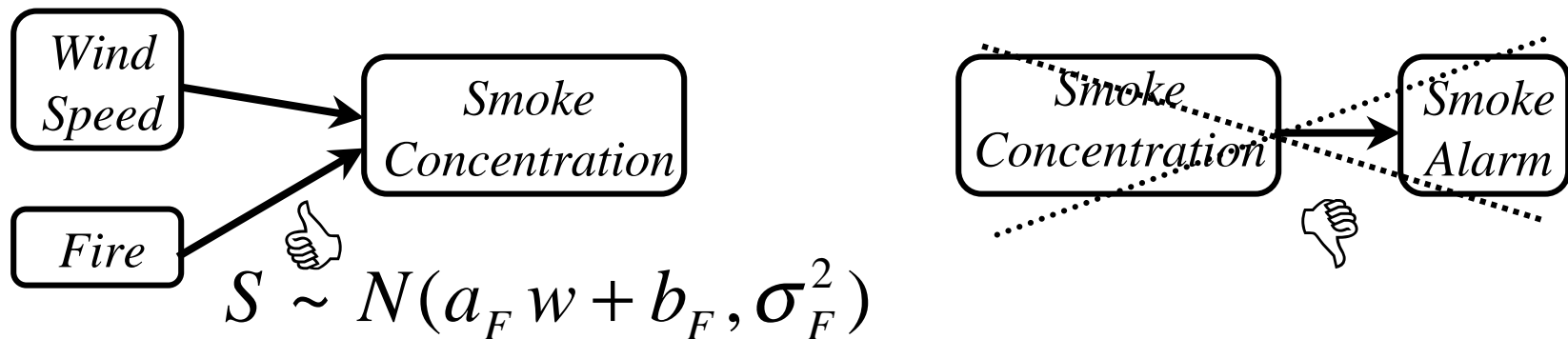    - discrete variables cannot depend on continuous

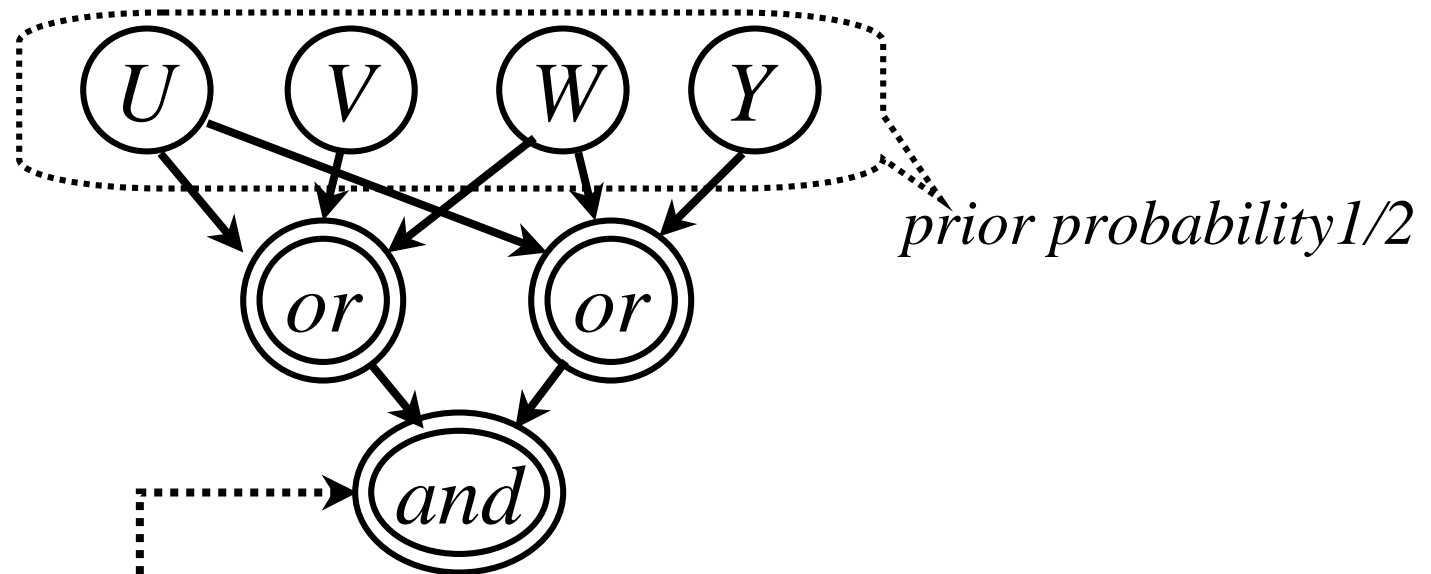

$$S \sim N(a_F w + b_F, \sigma_F^2)$$

- These techniques do not work for general hybrid networks.
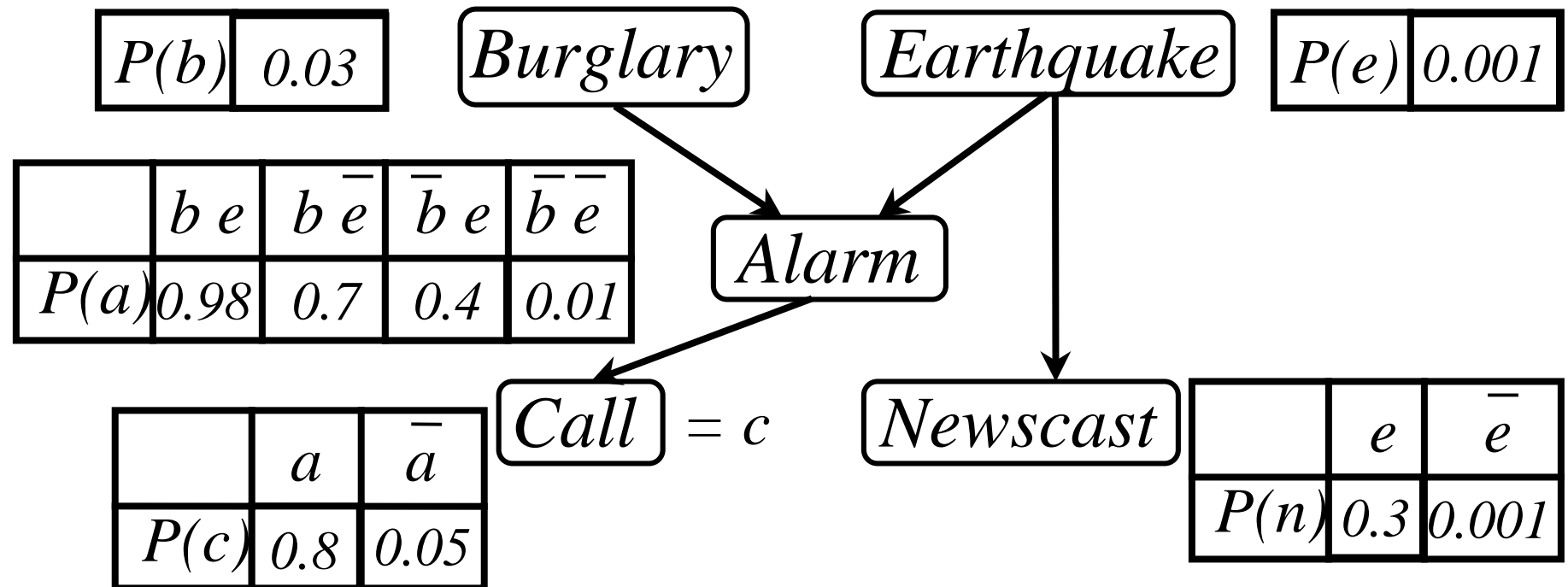
# Computational complexity

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

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



*prior probability 1/2*

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

72

# Stochastic simulation

| P(b) | 0.03 |
|------|------|

| Burglary |
|----------|

| Earthquake |
|------------|

| P(e) | 0.001 |
|------|-------|

| | $b\ e$ | $b\ \bar{e}$ | $\bar{b}\ e$ | $\bar{b}\ \bar{e}$ |
|------|------|------|------|------|
| P(a) | 0.98 | 0.7 | 0.4 | 0.01 |

| Alarm |
|-------|

| | $a$ | $\bar{a}$ |
|------|-----|-----|
| P(c) | 0.8 | 0.05 |

| Call | = c

| Newscast |
|----------|

| | $e$ | $\bar{e}$ |
|------|-----|-----|
| P(n) | 0.3 | 0.001 |

**Samples:**

| B | E | A | C | N |
|---|---|---|---|---|
| $\bar{b}$ | $e$ | $a$ | $c$ | $\bar{n}$ |
| ~~$b$~~ | ~~$\bar{e}$~~ | ~~$a$~~ | ~~$\bar{c}$~~ | ~~$n$~~ |
| ⋮ | | | | |

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

# Likelihood weighting

**Burglary**     **Earthquake**

**Alarm**

**Call** = c     **Newscast**

|  | $a$ | $\bar{a}$ |
|---|---|---|
| $P(c)$ | 0.8 | 0.05 |

Samples:

| $B$ | $E$ | $A$ | $C$ | $N$ | weight |
|---|---|---|---|---|---|
| $\bar{b}$ | $e$ | $a$ | $c$ | $\bar{n}$ | 0.8 |
| $b$ | $\bar{e}$ | $\bar{a}$ | $c$ | $n$ | 0.95 |
| ⋮ | | | | | |

$$P(b/c) = \frac{\text{weight of samples with } B=b}{\text{total weight of samples}}$$
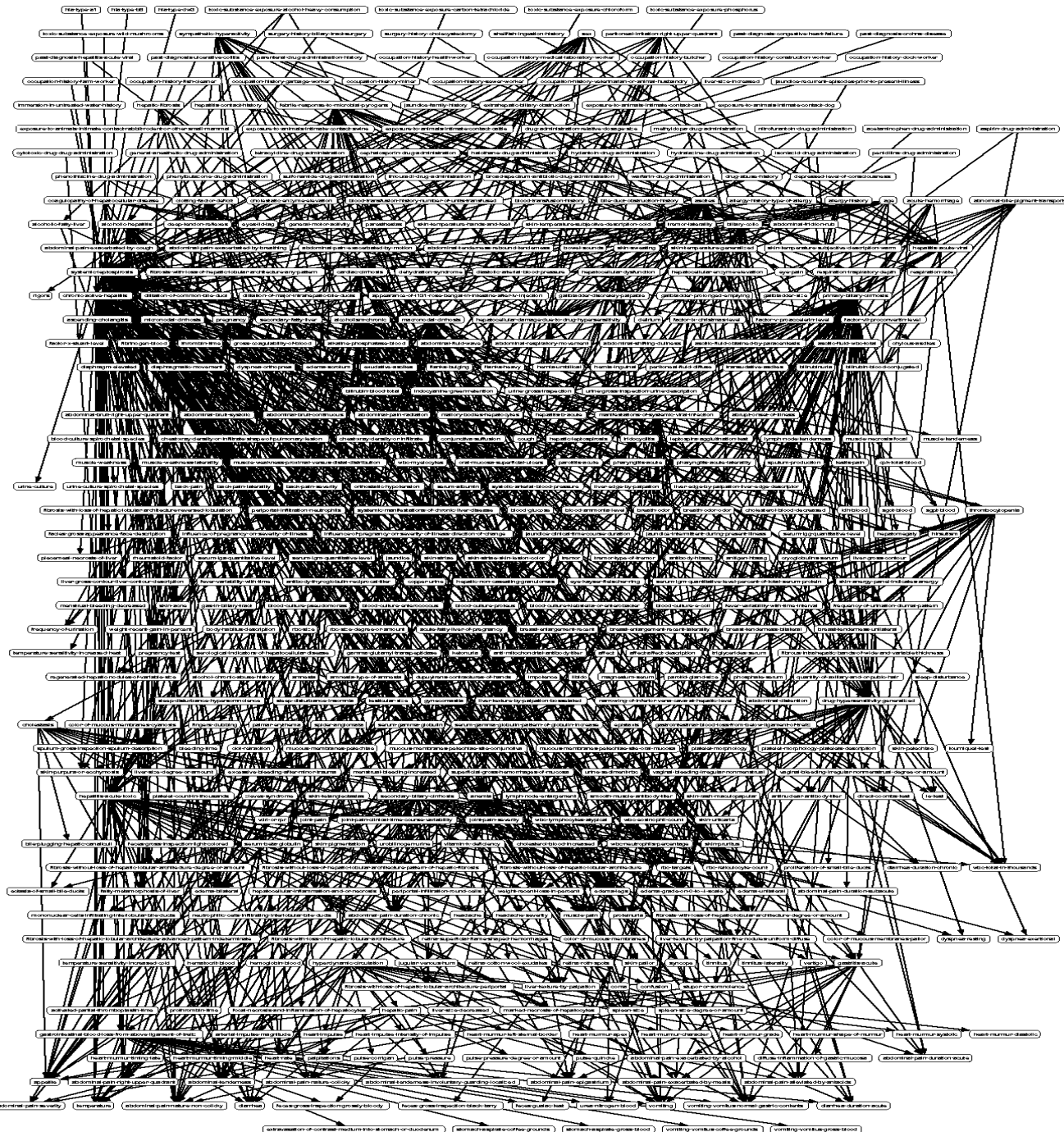
# Other approaches

■ **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

# CPCS
# Network

# Course Contents

■ Concepts in Probability

■ Bayesian Networks

■ Inference

» Decision making

■ Learning networks from data

■ Reasoning over time

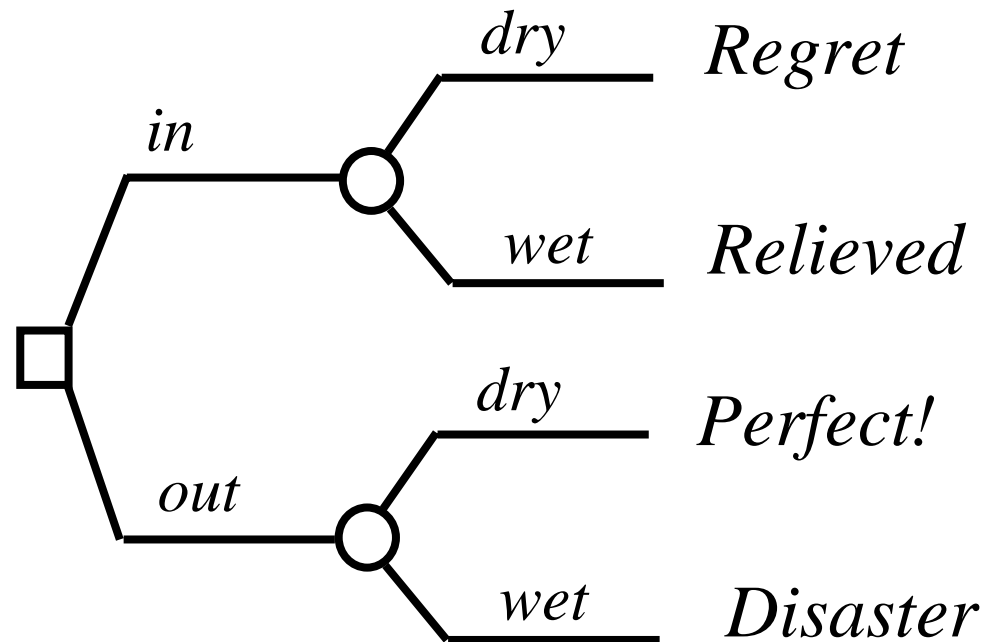■ Applications

# Decision making

■ Decisions, Preferences, and Utility functions

■ Influence diagrams

■ Value of information

# Decision making

■ Decision - an irrevocable allocation of domain resources

■ Decision should be made so as to maximize expected utility.

■ View decision making in terms of

◆ Beliefs/Uncertainties

◆ Alternatives/Decisions

◆ Objectives/Utilities

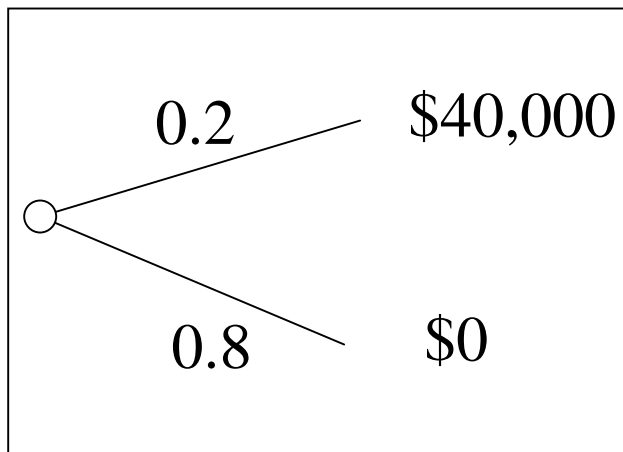# A Decision Problem

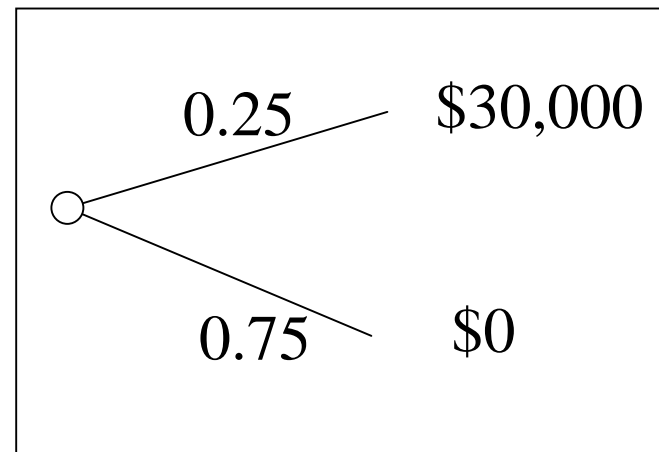Should I have my party inside or outside?

# Value Function

■ A numerical score over all possible states of the world.

| Location? | Weather? | Value |
|-----------|----------|-------|
| in        | dry      | $50   |
| in        | wet      | $60   |
| out       | dry      | $100  |
| out       | wet      | $0    |

# Preference for Lotteries



$$0.2 \quad \$40{,}000$$
$$0.8 \quad \$0$$

$$\succ$$
$$\succ$$
$$\approx$$

$$0.25 \quad \$30{,}000$$
$$0.75 \quad \$0$$

# Desired Properties for Preferences over Lotteries

If you prefer $100 to $0 and $p < q$ then



(always)

# Expected Utility

Properties of preference $\Rightarrow$
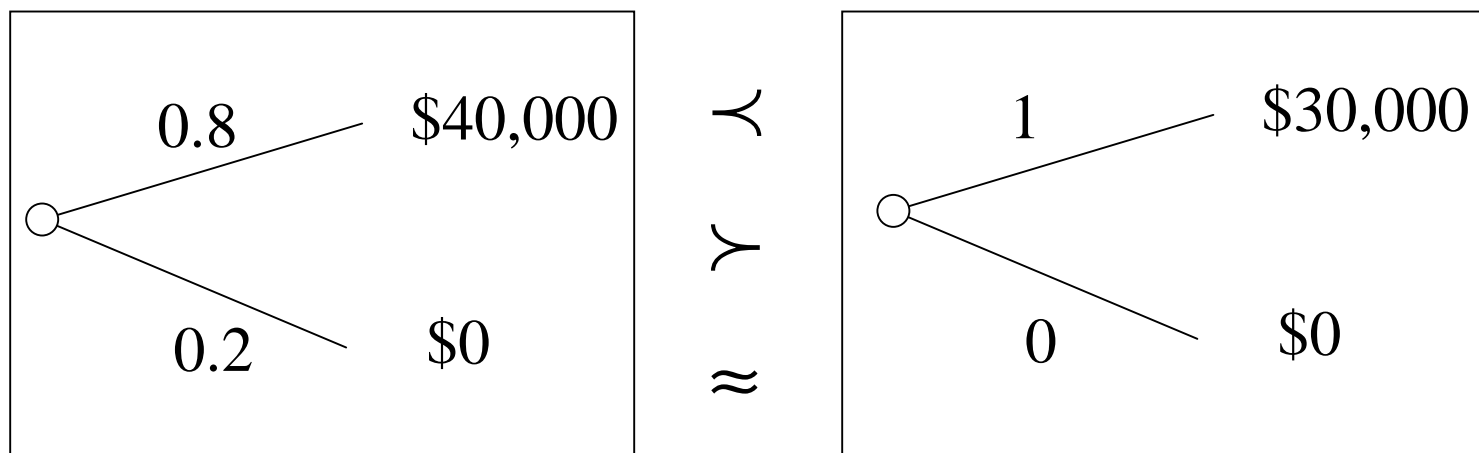     existence of function $U$, that satisfies:



$$\prec$$

iff

$$\Sigma_i \, p_i \, U(x_i) \quad < \quad \Sigma_i \, q_i \, U(y_i)$$

84

# Some properties of U



$$0.8 \quad \$40,000$$
$$0.2 \quad \$0$$

$$\succ$$
$$\succ$$
$$\approx$$

$$1 \quad \$30,000$$
$$0 \quad \$0$$

$$\Rightarrow U \neq \text{ monetary payoff}$$

# Attitudes towards risk

$U$

$U(\$500)$

$U(l)$

$l:$

.5     $1,000

.5     $0

0    400   500      1000    $ reward

Certain equivalent

insurance/risk premium

| U convex | risk averse |
|---|---|
| U concave | risk seeking |
| U linear | risk neutral |

# Are people rational?

0.2 — $40,000
0.8 — $0

$\succ$

0.25 — $30,000
0.75 — $0

$0.2 \cdot U(\$40k) \quad > \quad 0.25 \cdot U(\$30k)$
$0.8 \cdot U(\$40k) \quad > \quad U(\$30k)$

0.8 — $40,000
0.2 — $0

$\prec$

1 — $30,000
0 — $0

$0.8 \cdot U(\$40k) \quad < \quad U(\$30k)$

87

# Maximizing Expected Utility



*dry*   $U(\$50)=.632$

0.7

*in*
.652

0.3   *wet*   $U(\$60)=.699$

*dry*   $U(\$100)=.865$

0.7

*out*
.605

0.3   *wet*   $U(\$0) = 0$

choose the action that maximizes expected utility

$EU(in) = 0.7 \cdot .632 + 0.3 \cdot .699 = .652$

$\Longrightarrow$   Choose *in*

$EU(out) = 0.7 \cdot .865 + 0.3 \cdot 0 = .605$

# Multi-attribute utilities
## (or: Money isn't everything)

■ **Many aspects of an outcome combine to determine our preferences.**

  ◆ vacation planning: cost, flying time, beach quality, food quality, …

  ◆ medical decision making: risk of death (micromort), quality of life (QALY), cost of treatment, …

■ **For rational decision making, must combine all relevant factors into single utility function.**

# Influence Diagrams

90

# Decision Making with Influence Diagrams

Earthquake

Burglary

Alarm

Newcast

Call

| Call? | Go Home? |
|---|---|
| Neighbor Phoned | Yes |
| No Phone Call | No |

Go Home?

Goods Recovered

Miss Meeting

Big Sale

Utility

Expected Utility of this policy is 100

# Value-of-Information

■ What is it worth to get another piece of information?

■ What is the increase in (maximized) expected utility if I make a decision with an additional piece of information?

■ Additional information (if free) cannot make you worse off.

■ There is no value-of-information if you will not change your decision.

92

# Value-of-Information in an Influence Diagram



How much better can we do when this arc is here?

93

# Value-of-Information is the increase in Expected Utility

Earthquake

Burglary

Alarm

Newcast

Call

| Phonecall? | Newscast? | Go Home? |
|------------|-----------|----------|
| Yes | Quake | No |
| Yes | No Quake | Yes |
| No | Quake | No |
| No | No Quake | No |

*Go Home?*

*Goods Recovered*

*Miss Meeting*

*Big Sale*

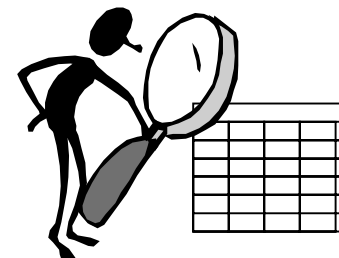*Utility*

Expected Utility of this policy is 112.5

# Course Contents

- Concepts in Probability

- Bayesian Networks

- Inference

- Decision making

» Learning networks from data

- Reasoning over time

- Applications

95

# Learning networks from data

- The learning task

- Parameter learning
  - Fully observable
  - Partially observable

- Structure learning

- Hidden variables

# The learning task

| B E A C N |
| --- |
| $\bar{b}$  e  $\bar{a}$  c  $\bar{n}$ |
| b  $\bar{e}$  $\bar{a}$  $\bar{c}$  n |
| ⋮ |

$$\Longrightarrow$$

Burglary    Earthquake

Alarm

Call    Newscast

*Input: training data*          *Output: BN modeling data*

■ Input: fully or partially observable data cases?

■ Output: parameters or also structure?

# Parameter learning: one variable

- **Unfamiliar coin:**
  - Let $\theta$ = bias of coin (long-run fraction of heads)

- **If $\theta$ known (given), then**
  - $P(X = heads \mid \theta) = \theta$

- **Different coin tosses independent given $\theta$**
  
  $\Rightarrow P(X_1, ..., X_n \mid \theta) = \theta^h (1-\theta)^t$
  
  $\underbrace{\phantom{P(X_1, ..., X_n)}}$
  
  $h$ heads, $t$ tails

# Maximum likelihood

- Input: a set of previous coin tosses
  - $X_1, \ldots, X_n = \{\underbrace{H, T, H, H, H, T, T, H, \ldots, H}_{h \text{ heads}, \, t \text{ tails}}\}$

- Goal: estimate $\theta$

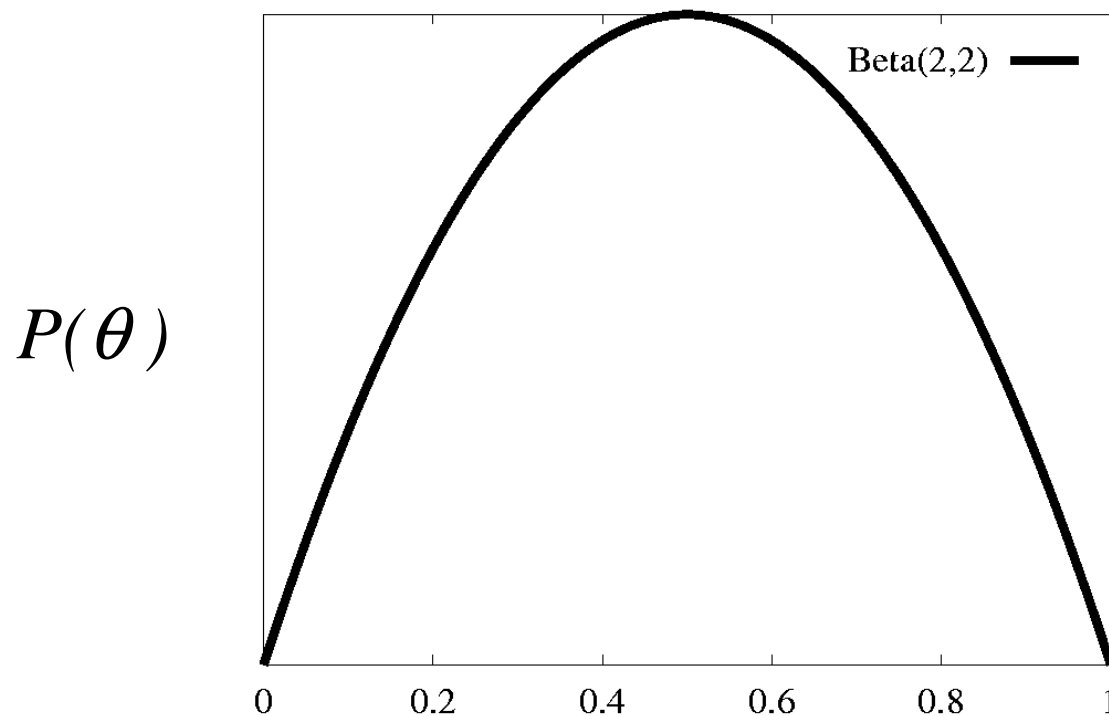- The likelihood $P(X_1, \ldots, X_n \,/\, \theta) = \theta^h \, (1-\theta)^t$

- The maximum likelihood solution is:

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

# Bayesian approach

Uncertainty about $\theta \Rightarrow$ distribution over its values



$P(\theta)$

Beta(2,2)

$\theta$

$$P(X = heads) = \int_{0}^{\infty} P(X = heads \mid \theta)P(\theta)d\theta = \int_{0}^{\infty} \theta\, P(\theta)\, d\theta$$

# Conditioning on data

$h$ heads, $t$ tails

$D$

$P(\theta)$ $\longrightarrow$ $P(\theta \mid D) \propto P(\theta) P(D \mid \theta)$

$$= P(\theta) \, \theta^h (1-\theta)^t$$

Beta(1,1)

*1 head*
*1 tail* $\longrightarrow$

Beta(2,2)

| 0 | 0.2 | 0.4 | 0.6 | 0.8 | 1 |

| 0 | 0.2 | 0.4 | 0.6 | 0.8 | 1 |

Good parameter distribution:

$$Beta(\alpha_h, \alpha_t) \propto$$

$$\theta^{\alpha_h - 1} (1 - \theta)^{\alpha_t - 1}$$

Beta(3,2) ——

Beta(10,10) ——

Beta(15,10) ——

* Dirichlet distribution generalizes Beta to non-binary variables.

# General parameter learning

■ A multi-variable BN is composed of several independent parameters ("coins").

$$A \longrightarrow B$$

*Three parameters:*
$$\theta_A, \ \theta_{B/a}, \ \theta_{B/\overline{a}}$$

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

*Max likelihood estimate of $\theta_{B/\overline{a}}$ would be:*

$$\theta^*_{B/\overline{a}} = \frac{\#data\ cases\ with\ b, \overline{a}}{\#data\ cases\ with\ \overline{a}}$$

# Partially observable data

| B | E | A | C | N |
|---|---|---|---|---|
| $\bar{b}$ | ? | a | c | ? |
| b | ? | $\bar{a}$ | ? | n |
| $\vdots$ | | | | |

Burglary → Alarm

Earthquake → Alarm

Earthquake → Newscast

Alarm → Call

■ Fill in missing data with "expected" value

◆ expected = distribution over possible values

◆ use "best guess" BN to estimate distribution

# Intuition

■ In fully observable case:

$$\theta^*_{n/e} = \frac{\#data\ cases\ with\ n,\ e}{\#data\ cases\ with\ e} = \frac{\Sigma_j I(n, e \mid d_j)}{\Sigma_j I(e \mid d_j)}$$

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

■ In partially observable case $I$ is unknown.

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

Problem: $\theta^*$ unknown.

# Expectation Maximization (EM)

Repeat :

- ■ Expectation (E) step
    - ◆ Use current parameters $\theta$ to estimate filled in data.

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

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

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

- ■ 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.

score

109

# 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:

  $$X \longrightarrow Y$$

**Max. likelihood network typically fully connected**

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

# Better scoring functions

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

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

■ 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.

# Hidden variables

■ There may be interesting variables that we never get to observe:

◆ topic of a document in information retrieval;

◆ user's current task in online help system.

■ Our learning algorithm should

◆ hypothesize the existence of such variables;

◆ learn an appropriate state space for them.

Randomly
scattered data

113

$E_1$

$E_3$

$E_2$

Actual data

# Bayesian clustering (Autoclass)

naïve Bayes model:



- ■ (hypothetical) class variable never observed
- ■ if we know that there are $k$ classes, just run EM
- ■ learned classes = clusters
- ■ Bayesian analysis allows us to choose $k$, trade off fit to data with model complexity

E₁

E₃

E₂

Clustered
distributions

116

# Detecting hidden variables

- Unexpected correlations ⇨ hidden variables.

Hypothesized model

Data model

"Correct" model

# Course Contents

- Concepts in Probability

- Bayesian Networks

- Inference

- Decision making

- Learning networks from data

  » Reasoning over time

- Applications

# Reasoning over time

■ Dynamic Bayesian networks

■ Hidden Markov models

■ Decision-theoretic planning

◆ Markov decision problems

◆ Structured representation of actions

◆ The qualification problem & the frame problem

◆ Causality (and the frame problem revisited)

# Dynamic environments

$State(t) \longrightarrow State(t+1) \longrightarrow State(t+2)$

- **Markov property:**
  - past independent of future given current state;
  - a conditional independence assumption;
  - implied by fact that there are no arcs $t \to t+2$.

# Dynamic Bayesian networks

■ State described via random variables.

■ Each variable depends only on few others.

# Hidden Markov model

■ An HMM is a simple model for a partially observable stochastic domain.



*State(t)* → *State(t+1)* — State transition model

*State(t)* → *Obs(t)*

*State(t+1)* → *Obs(t+1)* — Observation model

# Hidden Markov models (HMMs)

## Partially observable stochastic environment:

- ■ Mobile robots:
  - ◆ states = location
  - ◆ observations = sensor input

- ■ Speech recognition:
  - ◆ states = phonemes
  - ◆ observations = acoustic signal

- ■ Biological sequencing:
  - ◆ states = protein structure
  - ◆ observations = amino acids

*0.15*

*0.05*

*0.8*

# HMMs and DBNs

■ HMMs are just very simple DBNs.

■ Standard inference & learning algorithms for HMMs are instances of DBN algorithms

   ◆ Forward-backward = polytree

   ◆ Baum-Welch = EM

   ◆ Viterbi = most probable explanation.

# Acting under uncertainty

## Markov Decision Problem (MDP)

agent observes state ⋯⋯⋯⋯⋯⋯⋯⋯

action model ⋯⋯⋯⋯⋯⋯⋯⋯

*Action(t)*    *Action(t+1)*

*State(t)*    *State(t+1)*    *State(t+2)*

*Reward(t)*    *Reward(t+1)*

■ Overall utility = sum of momentary rewards.

■ Allows rich preference model, e.g.:

$$\text{rewards corresponding to "get to goal asap"} = \begin{cases} +100 & \text{goal states} \\ -1 & \text{other states} \end{cases}$$

125

# Partially observable MDPs

agent observes
***Obs***, not state

***Obs*** depends
on state

*Action(t)* → *Action(t+1)*

*Obs(t)*          *Obs(t+1)*

*State(t)* → *State(t+1)* → *State(t+2)*

*Reward(t)*          *Reward(t+1)*

- ■ The optimal action at time t depends on the entire history of previous observations.

- ■ Instead, a distribution over *State(t)* suffices.

# Structured representation

Preconditions ···

Move:

Position(t) → Position(t+1)

Direction(t) ·······> Direction(t+1) ······· Effects

Holding(t) ·······> Holding(t+1)

Turn:

Position(t) ·······> Position(t+1)

Direction(t) → Direction(t+1)

Holding(t) ·······> Holding(t+1)

## Probabilistic action model
- allows for exceptions & qualifications;
- persistence arcs: a solution to the frame problem.

127

# Causality

■ Modeling the effects of interventions

■ Observing vs. "setting" a variable

■ A form of persistence modeling

# Causal Theory

```
   ┌─────────────┐
   │ Temperature │
   └─────────────┘
          │
          ▼
 ┌───────────────────┐
 │  Distributor Cap  │
 └───────────────────┘
          │
          ▼
   ┌─────────────┐
   │  Car Starts │
   └─────────────┘
```

Cold temperatures can cause the distributor cap to become cracked.

If the distributor cap is cracked, then the car is less likely to start.

# Setting vs. Observing

Temperature

Distributor Cap

Car Starts

The car does not start.
Will it start if we
replace the distributor?

# Predicting the effects of interventions

```
   ( Temperature )
         |
         v
 ( Distributor Cap )
         |
         v
   ( Car Starts )
```

The car does not start. Will it start if we replace the distributor?

What is the probability that the car will start if I replace the distributor cap?

# Mechanism Nodes



| M$_{start}$ | Distributor | Starts? |
|---|---|---|
| Always Starts | Cracked | Yes |
| Always Starts | Normal | Yes |
| Never Starts | Cracked | No |
| Never Starts | Normal | No |
| Normal | Cracked | No |
| Normal | Normal | Yes |
| Inverse | Cracked | Yes |
| Inverse | Normal | No |

# Persistence

Pre-action

Post-action

**Temperature**

**Temperature**

**Dist**

$M_{start}$

**Start**

Persistence
arc

Observed
Abnormal

$M_{start}$

**Dist**

Set to
Normal

**Start**

Assumption:The mechanism relating *Dist* to *Start* is
unchanged by replacing the *Distributor*.

# 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.

137

# Studies of Pathfinder Diagnostic Performance

■ Naïve Bayes performed considerably better than certainty factors and Dempster-Shafer Belief Functions.

■ Incorrect zero probabilities caused 10% of cases to be misdiagnosed.

■ Full Bayesian network model with feature dependencies did best.

138

# Commercial system: Integration

- Expert System with advanced diagnostic capabilities
  - uses key features to form the differential diagnosis
  - recommends additional features to narrow the differential diagnosis
  - recommends features needed to confirm the diagnosis
  - explains correct and incorrect decisions
- Video atlases and text organized by organ system
- "Carousel Mode" to build customized lectures
- Anatomic Pathology Information System

# 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

# Single Fault approximation

142

# On Parenting: Selecting problem

**Describe the child** in the drop-down boxes at the right. Relevant information will appear below.

**Age:** Toddler ▼   **Sex:** Female ▼

**Complaint:** [                    ] ▼

| |
|---|
| Abdominal pain ▲ |
| Abnormal control of body movements |
| Biting or hitting |
| Blood in stool |
| Blood in urine |
| Blood in vomit |
| Bluish or purplish skin |
| Breath-holding |
| Breathlessness or difficulty breathing |
| Colic or gas pain |
| Constipation |
| Cough |
| Delayed development |
| Delayed speech |
| Diarrhea |
| Difficulty swallowing ▼ |

This feature is designed to help you find information relevant to questions you answer about childhood symptoms. Keep in mind that any information you find is in no way comprehensive.

Also, this feature is NOT intended to be used to diagnose medical conditions or replace the advice of a healthcare professional. Always contact your healthcare provider for medical advice.

143

# Performing diagnosis/indexing

**Describe the child** in the drop-down boxes at the right. Relevant information will appear below.

**Age:** [Toddler ▼]  **Sex:** [Female ▼]

**Complaint:** [Abdominal pain ▼]

Localized pain: Can the child localize, or point to, the site of the pain?
- ○ No, unable to localize
- ○ Below the navel to the child's left
- ○ Above the child's navel
- ○ Either of the child's sides
- ○ Below the navel to the child's right
- ○ Above the navel to the child's right
- ○ Above the navel to the child's left
- ○ Don't Know

[ Start Over ]   [ Review ]

[ Next>> ]   [ Finish ]

## Results so far

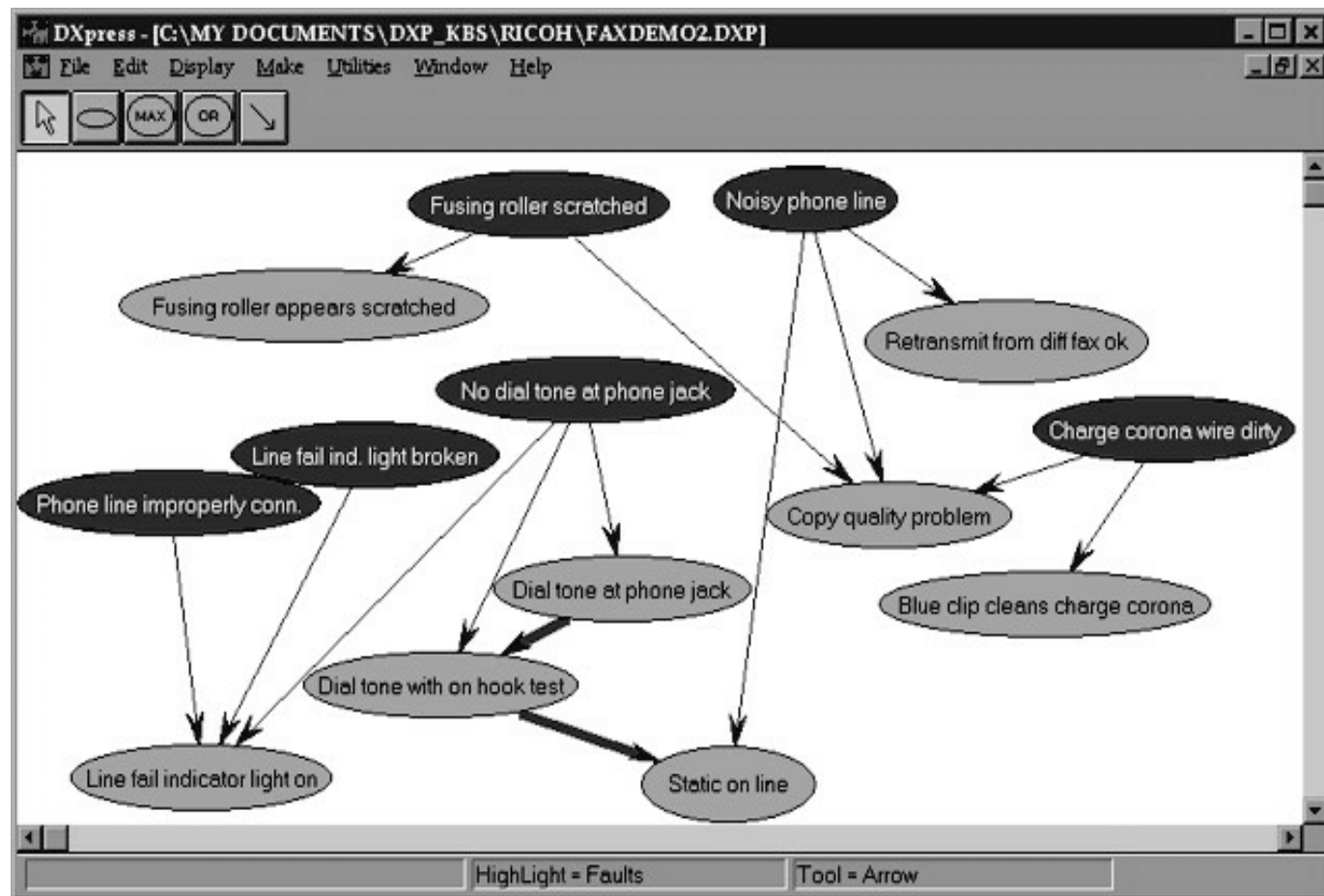| Disorder | Relevance |
|----------|-----------|
| Viral gastroenteritis | ▓▓▓▓░░ |
| Psychosomatic pain | ▓▓▓░░░ |
| Urinary tract infection | ▓░░░░░ |
| Other | ▓░░░░░ |

# RICOH Fixit

■ Diagnostics and information retrieval

145

# FIXIT: Ricoh copy machine

146

# Online Troubleshooters



Microsoft Technical Support Troubleshooters - Microsoft Internet Explorer

File  Edit  View  Go  Favorites  Help

Back  Forward  Stop  Refresh  Home  Search  Favorites  Print  Font  Mail  Edit

Address http://www.microsoft.com/support/tshooters.htm

Links  Computer News  Sidewalk  Sportszone

MICROSOFT | PRODUCTS | SEARCH | SUPPORT | SHOP | WRITE US  Microsoft

Microsoft
**Technical Support**

Site Map
Guided Tour with the Support Wizard
Knowledge Base
Troubleshooting Wizards
Frequently Asked Questions
Help Files, Service Packs, & Other Files
Newsgroups
Support Options & Phone Numbers
Submit a Question to a Support Engineer
Support Highlights

Best experienced with
Microsoft Internet Explorer
FREE
Click here to start.

## Troubleshooting Wizards

Microsoft now offers advanced inference engine technology to help you easily troubleshoot problems with Microsoft products. These Troubleshooting Wizards are the electronic version of our best engineers. Try them and see!

Choose a Troubleshooting Wizard from the list below, then click *Next*.

○ **Access:** Conversion Troubleshooter
○ **Access:** Help exporting to Active Server Pages and viewing them in your web browser.
○ **DirectX:** Help with common issues when you run applications with DirectX
○ **Excel:** PivotTable Troubleshooter
○ **Excel:** Workbook Troubleshooter
○ **Excel:** Video Troubleshooter
○ **Exchange:** Directory Synchronization Troubleshooter
○ **Exchange:** Internet Mail Connector Troubleshooter
○ **Games:** Help with problems when you install or start your Games
○ **Games:** Help with games display problems
○ **Games:** Help with sound problems when you run your games
○ **Games:** Help with common issues when you run your games
○ **Internet Explorer 3.0:** Help browsing the web

147

# Define Problem

## Troubleshooting Wizards

## Print Troubleshooter

The Print Troubleshooter lists recommended troubleshooting steps in the order of great
benefit and least cost to you (the user).

---

### What type of problem are you having?

- ⦿ My document didn't print at all.
- ○ Graphics look incomplete or incorrect.
- ○ Fonts are missing or do not look as they did on the screen.
- ○ The printout is garbled or contains garbage.
- ○ I only got part of the page I expected.
- ○ Printing is unusually slow.

[ Next ]

at www.microsoft.com

# Gather Information

**Troubleshooting Wizards**

## Print Troubleshooter

This table tracks your status in the troubleshooting process. If you need to change you
to a question, you can do so below:

| Problem: | Print Output |
| --- | --- |

**Are you printing from an MS-DOS-based or a Windows-based application?**

- ○ I am printing from MS-DOS or from an MS-DOS application.
- ● I am printing from a Windows application.
- ○ I don't want to do this now.

[ Next ]

149

# Get Recommendations

## Print Troubleshooter

This table tracks your status in the troubleshooting process. If you need to change yo[u]
to a question, you can do so below:

| Problem: | Print Output |
|---|---|
| Print Environment: | ○ MS-DOS ⦿ Windows ○ Unknown |
| Printing over Network: | ○ No (Local printer) ⦿ Yes (Network printer) ○ Un |
| Printer Driver Set Offline: | ⦿ Online ○ Unknown |

**Is your printer turned on and on-line?**

1. Make sure the printer is properly plugged into a power outlet.
2. Turn on the printer's power switch.
3. Make sure the printer is **on line**. Most printers have an On Line button with a li[ght]
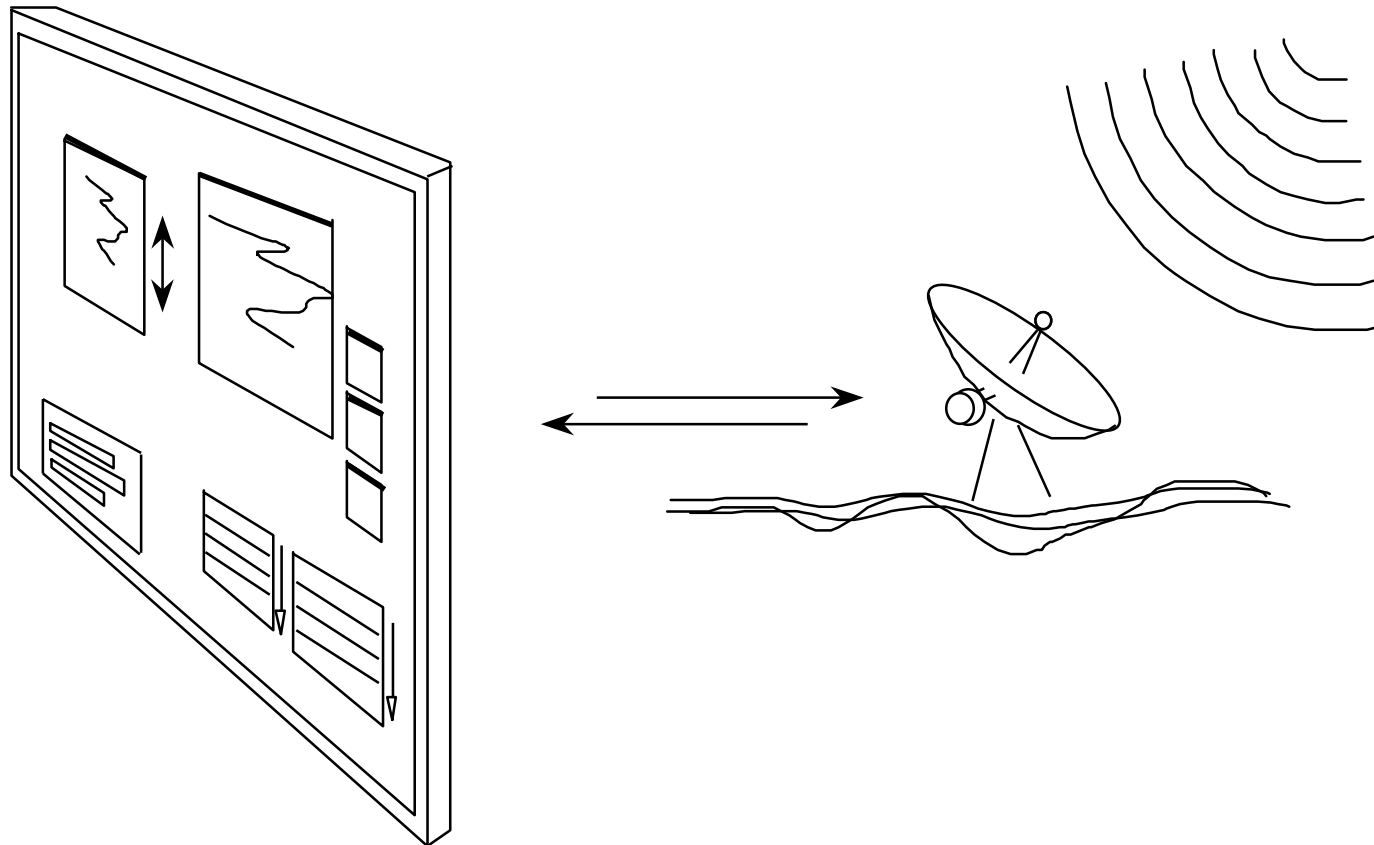4. Make sure the light is on.

If you need more information on any of these steps, consult your printer's manual.

- ○ It worked! I turned it on and now I can print.
- ○ Yes, my printer is on, but it still won't print.
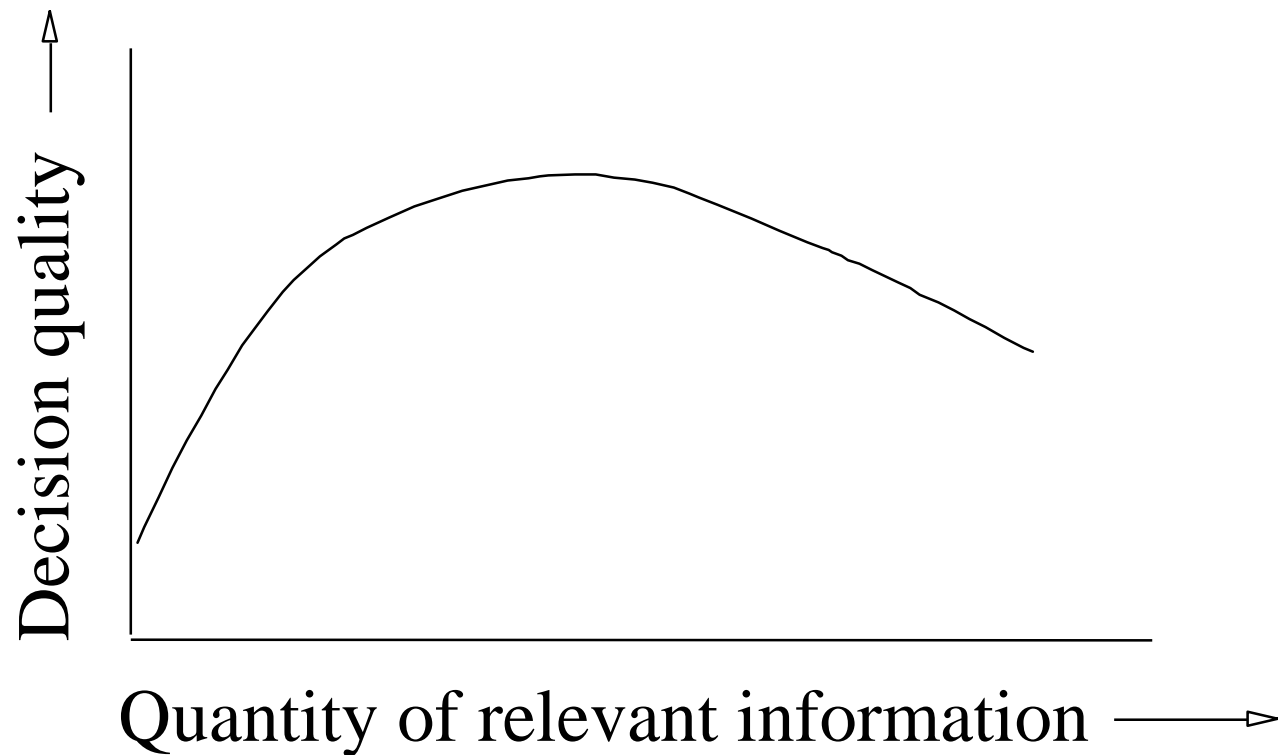- ○ I don't want to do this now.

Next

# Vista Project: NASA Mission Control

Decision-theoretic methods for display for high-stakes aerospace
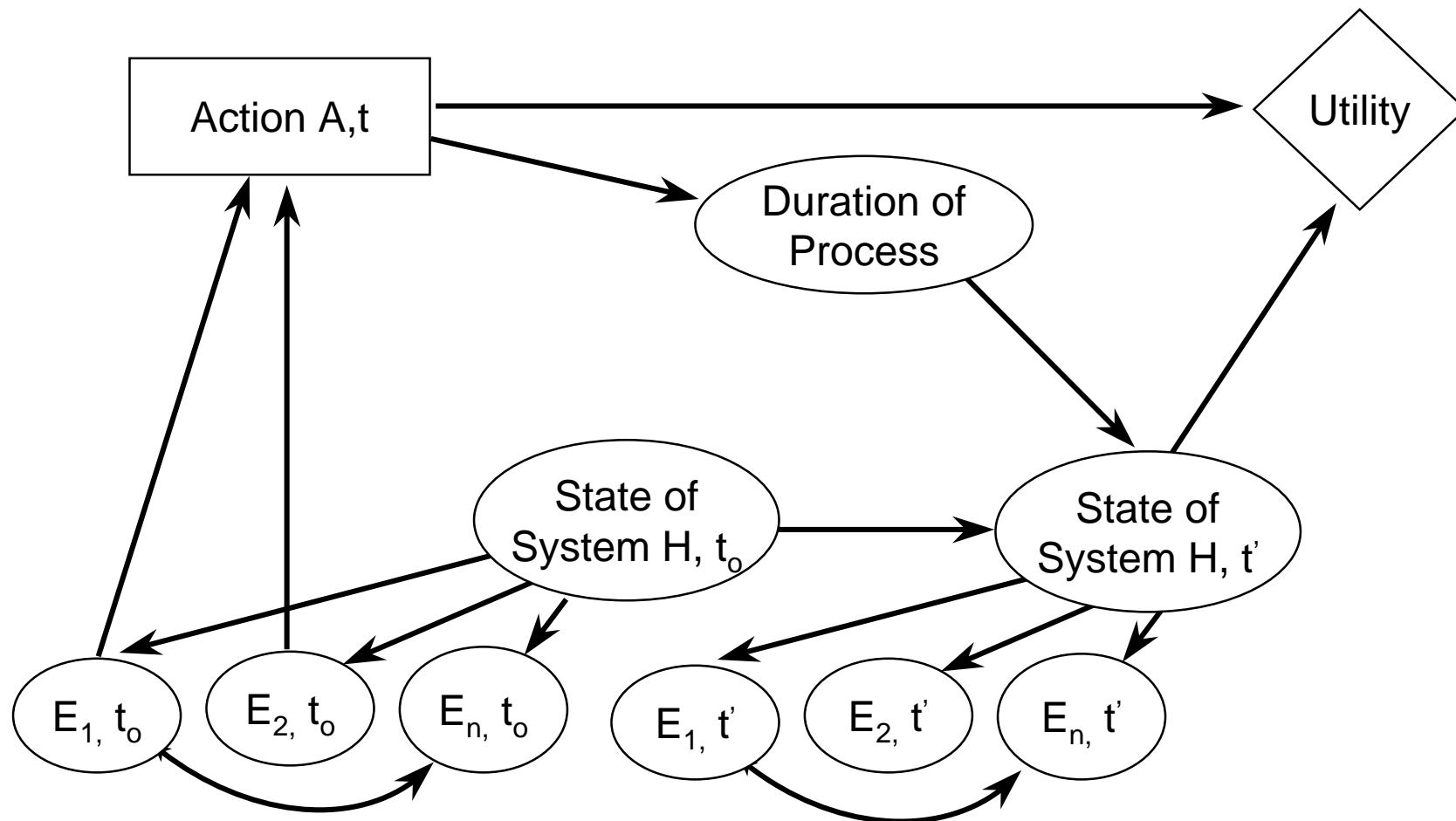  decisions

# Costs & Benefits of Viewing Information



Decision quality ⟶

Quantity of relevant information ⟶

# Status Quo at Mission Control
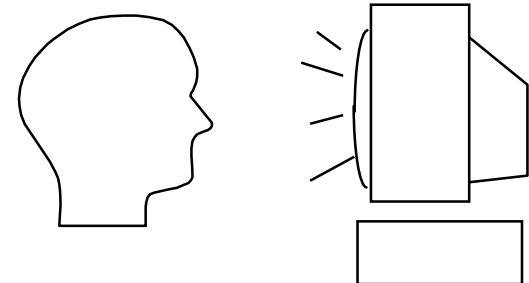
# Time-Critical Decision Making

- **Consideration of time delay in temporal process**

# Simplification: Highlighting Decisions

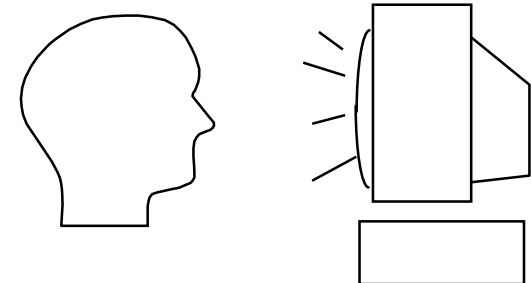■ Variable threshold to control amount of highlighted information

|  |  |  |
|---:|---|---|
| Oxygen | 15.6 | 14.2 |
| Fuel Pres | 10.5 | 11.8 |
| Chamb Pres | 5.4 | 4.8 |
| He Pres | 17.7 | 14.7 |
| Delta v | 33.3 | 63.3 |
|  |  |  |
| Oxygen | 10.2 | 10.6 |
| Fuel Pres | 12.8 | 12.5 |
| Chamb Pres | 0.0 | 0.0 |
| He Pres | 15.8 | 15.7 |
| Delta v | 32.3 | 63.3 |

# Simplification: Highlighting Decisions

■ Variable threshold to control amount of highlighted information

| | | |
|---|---|---|
| Oxygen | 15.6 | 14.2 |
| Fuel Pres | 10.5 | 11.8 |
| Chamb Pres | 5.4 | 4.8 |
| He Pres | 17.7 | 14.7 |
| Delta v | 33.3 | 63.3 |
| | | |
| Oxygen | 10.2 | 10.6 |
| Fuel Pres | 12.8 | 12.5 |
| Chamb Pres | 0.0 | 0.0 |
| He Pres | 15.8 | 15.7 |
| Delta v | 32.3 | 63.3 |

# Simplification: Highlighting Decisions

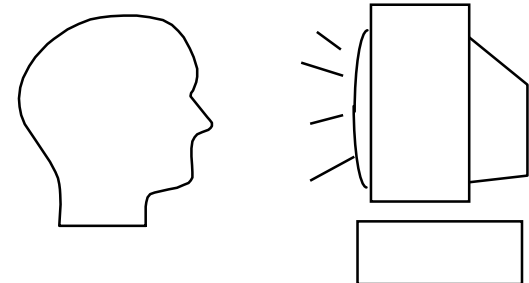■ Variable threshold to control amount of highlighted information

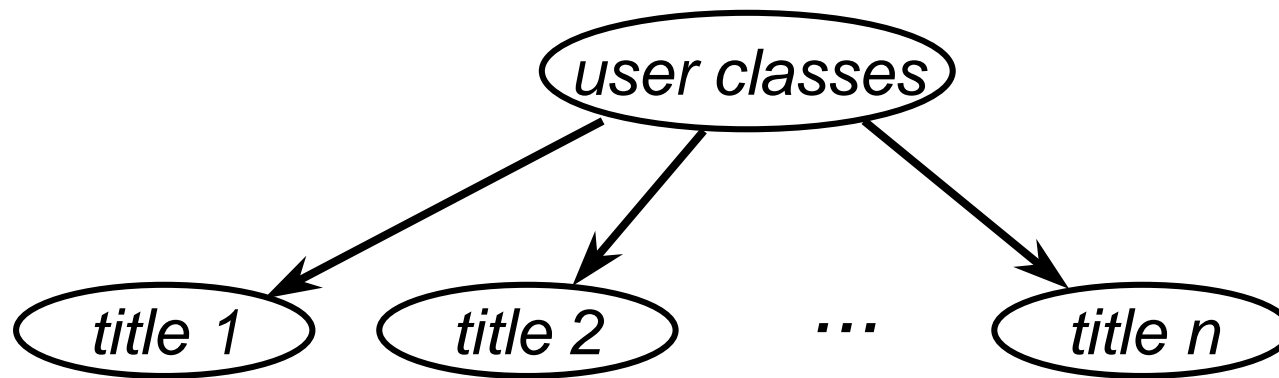| | | |
|---|---|---|
| Oxygen | 15.6 | 14.2 |
| Fuel Pres | 10.5 | 11.8 |
| Chamb Pres | 5.4 | 4.8 |
| He Pres | 17.7 | 14.7 |
| Delta v | 33.3 | 63.3 |
| | | |
| Oxygen | 10.2 | 10.6 |
| Fuel Pres | 12.8 | 12.5 |
| Chamb Pres | 0.0 | 0.0 |
| He Pres | 15.8 | 15.7 |
| Delta v | 32.3 | 63.3 |

# 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(Like title i | Like title j, Like title k)*

# Applying Bayesian clustering



|  | class1 | class2 | ... |
|---|---|---|---|
| title1 | p(like)=0.2 | p(like)=0.8 | |
| title2 | p(like)=0.7 | p(like)=0.1 | |
| title3 | p(like)=0.99 | p(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

## 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?

## Readers of "Softer" News (12%):

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

161

# 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
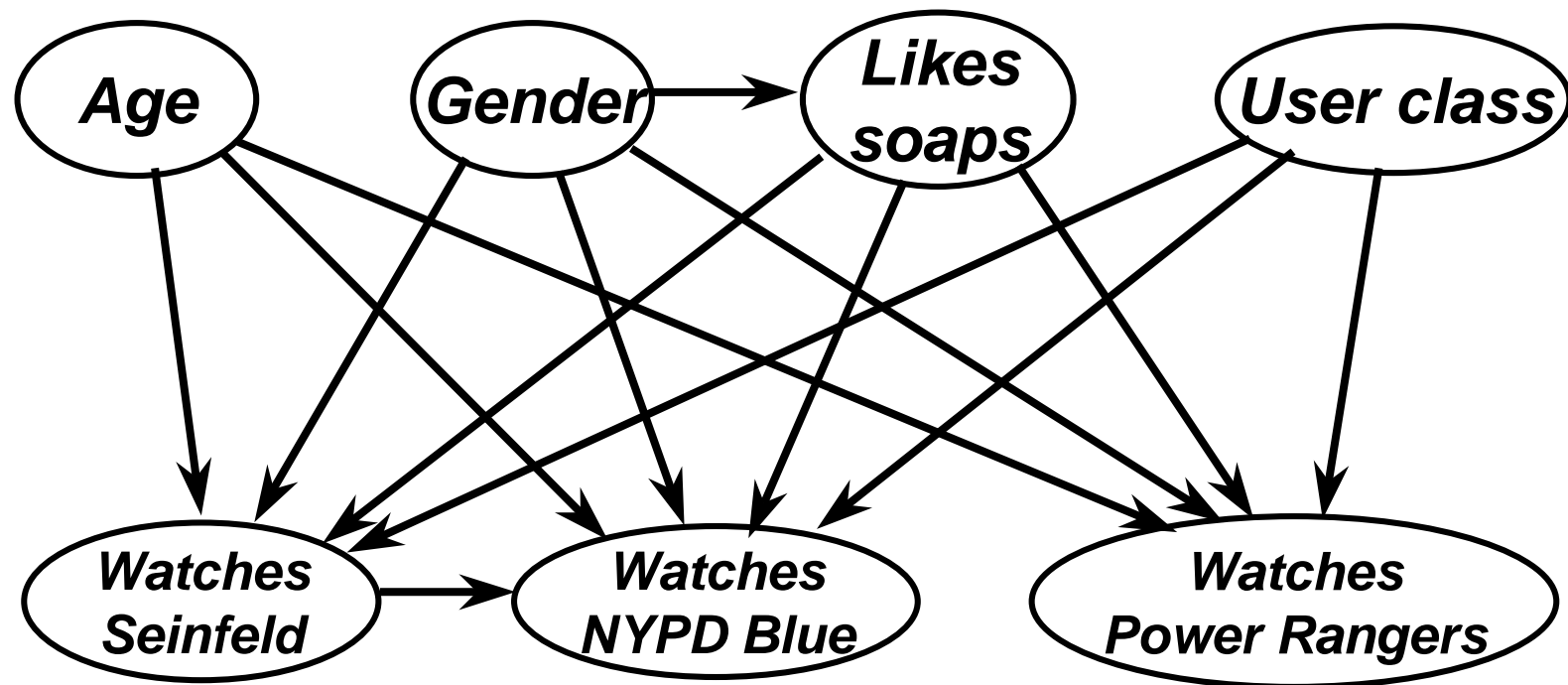- NBC nightly news
- Later with Kinnear
- Seinfeld

**Class 4**
- 60 minutes
- NBC nightly news
- CBS eve news
- Murder she wrote
- Matlock

**Class 5**
- Seinfeld
- Friends
- Mad about you
- ER
- Frasier

162

# Richer model



163

# 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.);
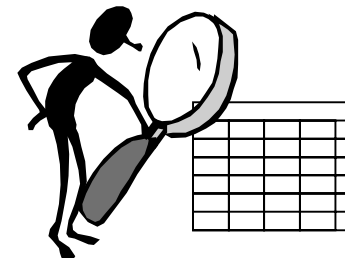  - parameter estimation from data.

# What's new?

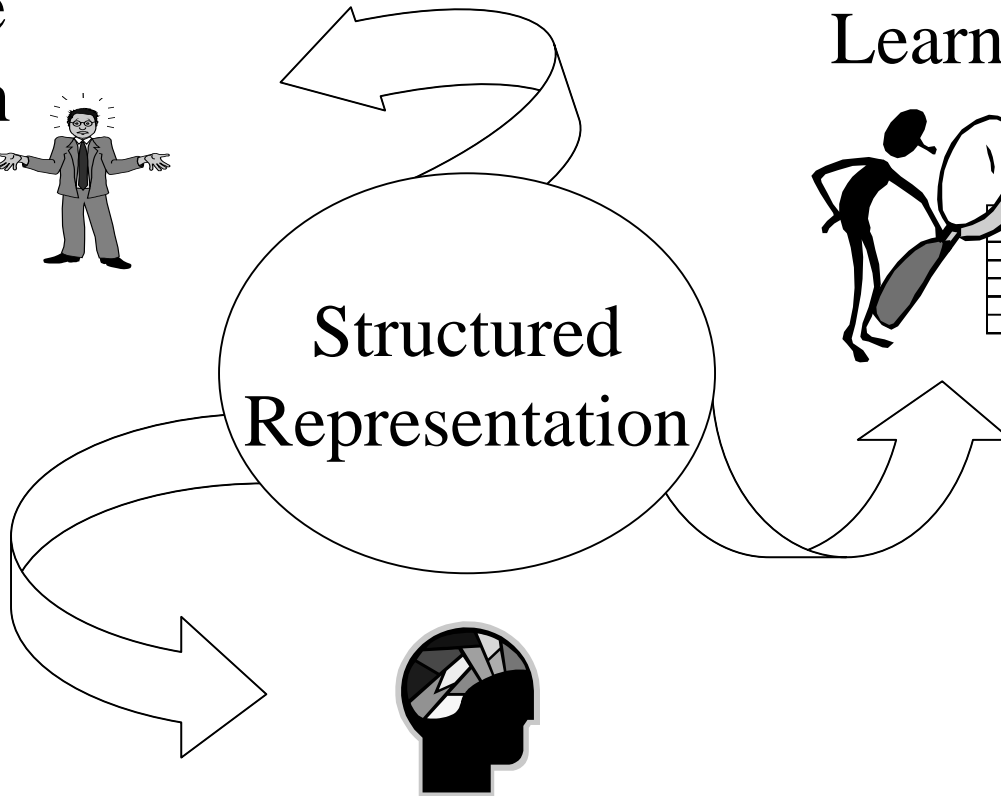Bayesian networks exploit domain structure to allow compact representations of complex models.

Knowledge
Acquisition

Learning

Structured
Representation

Inference

165

# Some Important AI Contributions

- Key technology for diagnosis.
- Better more coherent expert systems.
- New approach to planning & action modeling:
  - planning using Markov decision problems;
  - new framework for reinforcement learning;
  - probabilistic solution to frame & qualification problems.
- New techniques for learning models from data.

166

# 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.

Structured Representation