

Introduction to Artificial Intelligence

Bayesian Networks

Jianmin Li

Department of Computer Science and Technology
Tsinghua University

Spring, 2024

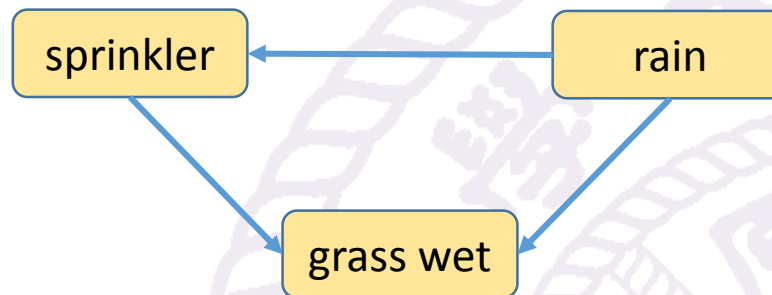
Outline

- Syntax
- Semantics
- Parameterized distributions
- Exact inference



Bayesian networks

- A simple, graphical notation for conditional independence assertions
- Compact specification of full joint distributions

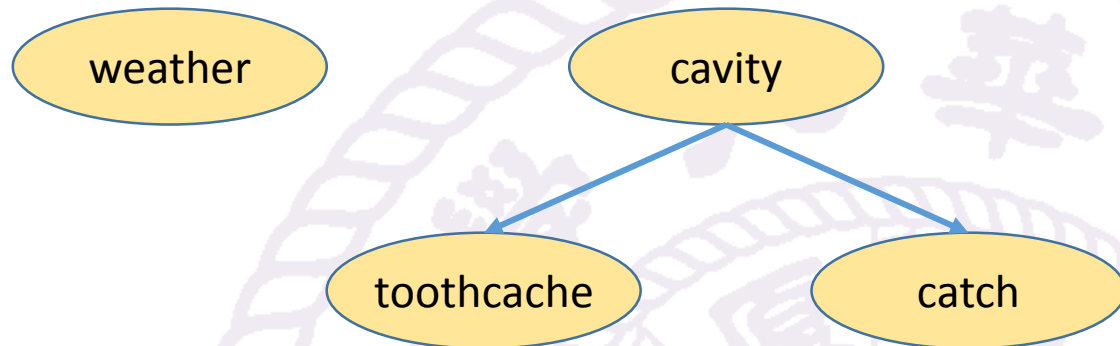


- Syntax:
 - a set of nodes, one per variable
 - a directed, acyclic graph (link “directly influences”)
 - a conditional distribution for each node given its parents

$$P(X_i | Parents(X_i))$$

Topology

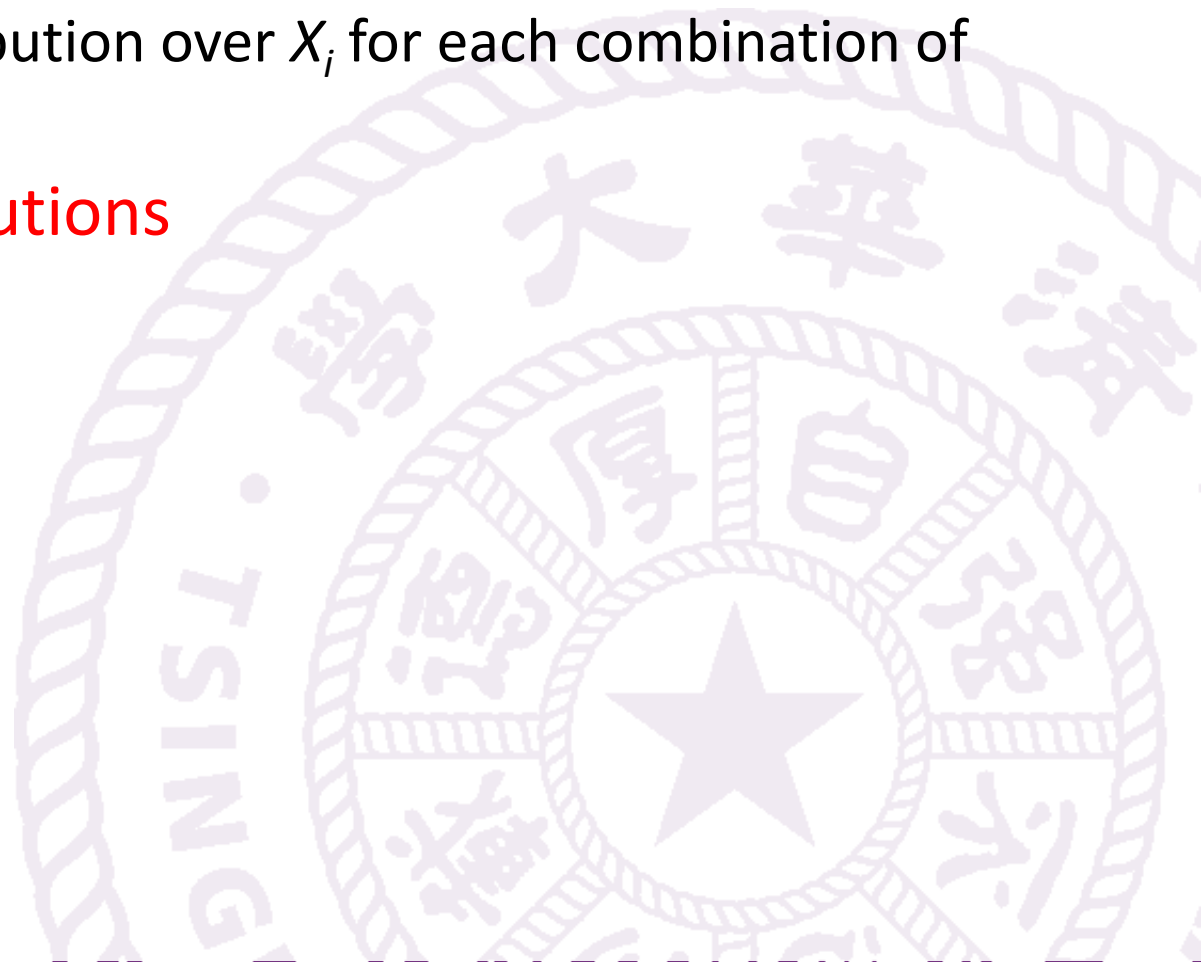
- Topology of network encodes conditional independence assertions



- *Weather* is independent of the other variables
- *Toothache* and *Catch* are conditionally independent given *Cavity*

Conditional distribution

- **Conditional probability table** (CPT)
 - giving the distribution over X_i for each combination of parent values
- **Canonical distributions**

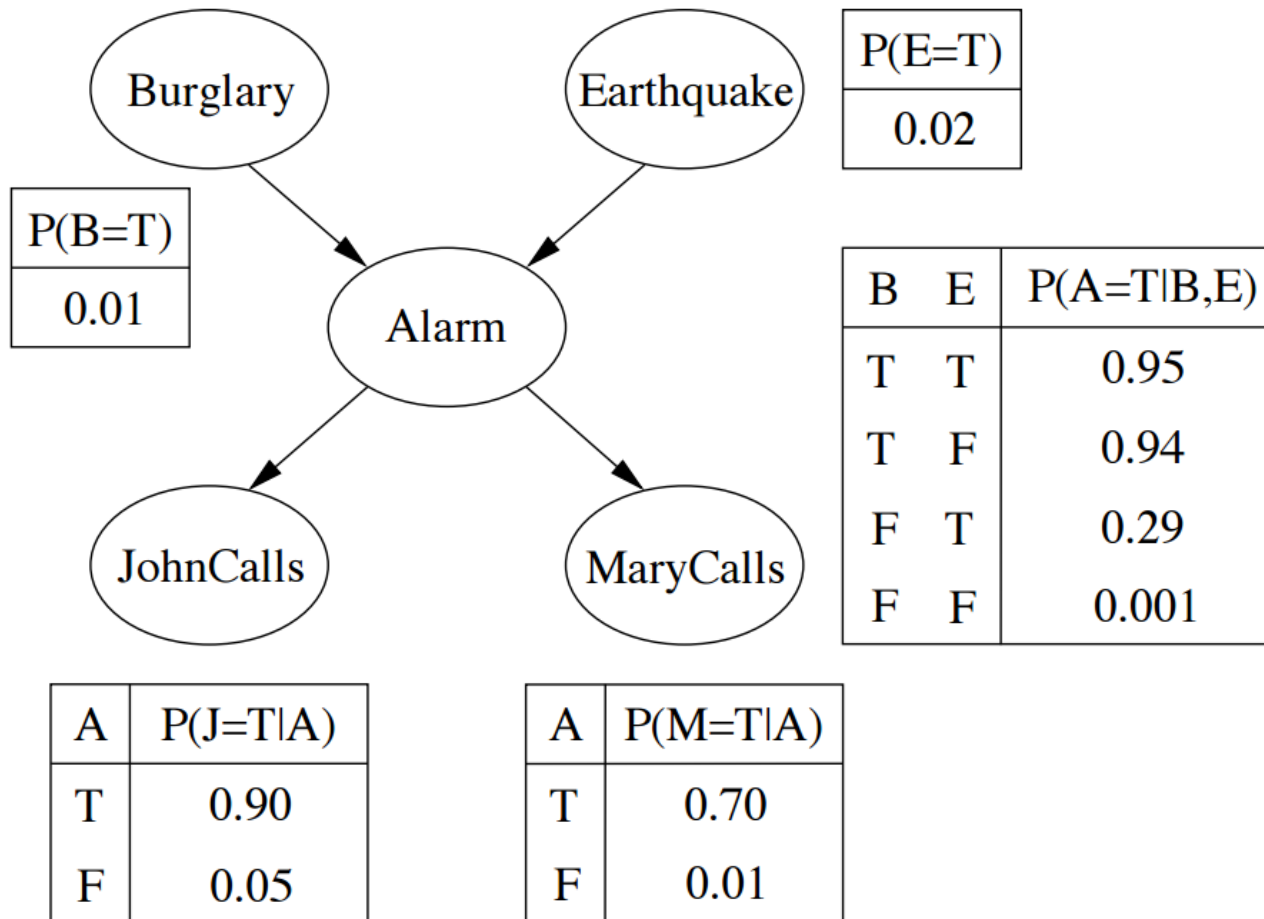


Example

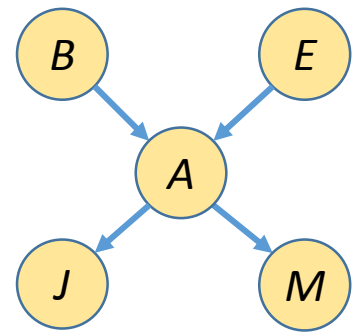


- I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?
- Variables: *Burglary*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*
- Network topology reflects "causal" knowledge:
 - A burglar can set the alarm off
 - An earthquake can set the alarm off
 - The alarm can cause Mary to call
 - The alarm can cause John to call

Example

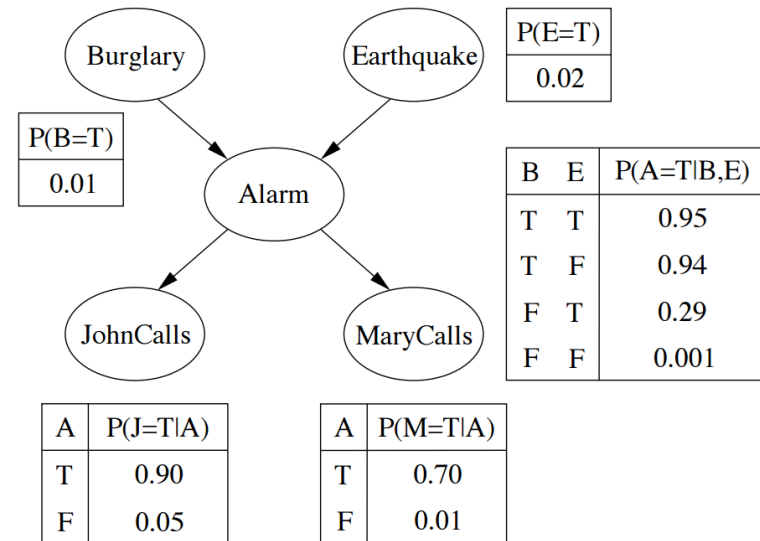


Compactness



- A CPT for Boolean X_i with k Boolean parents has 2^k rows for the combinations of parent values
- Each row requires one number p for $X_i = \text{true}$ (the number for $X_i = \text{false}$ is just $1-p$)
- If each variable has no more than k parents, the complete network requires $O(n \cdot 2^k)$ numbers
- I.e., grows linearly with n , vs. $O(2^n)$ for the full joint distribution
- For burglary net, $1+1+4+2+2=10$ numbers (vs. $2^5-1=31$)

Global semantics



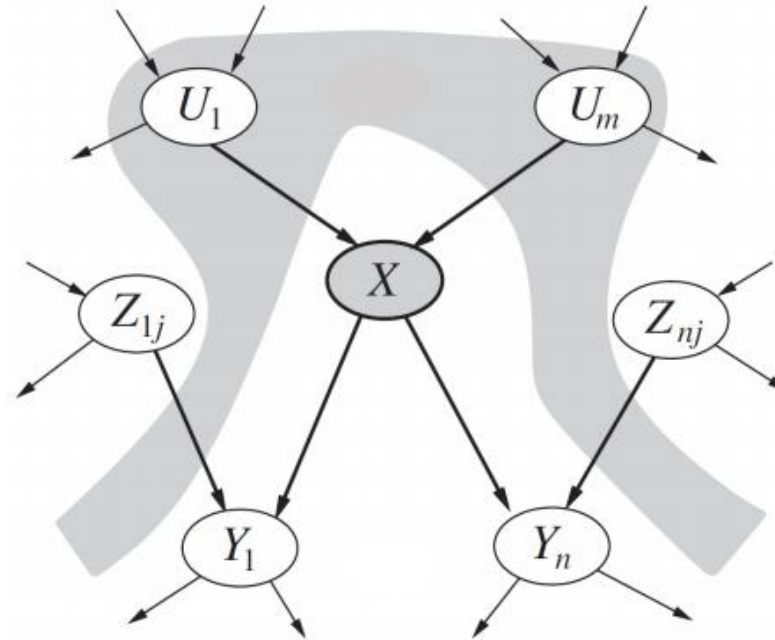
- Defines the full joint distribution as the product of the local conditional distributions

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

- e.g

$$\begin{aligned}
 & P(j \wedge m \wedge a \wedge \neg b \wedge \neg e) \\
 = & P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) \\
 = & 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\
 \approx & 0.000628
 \end{aligned}$$

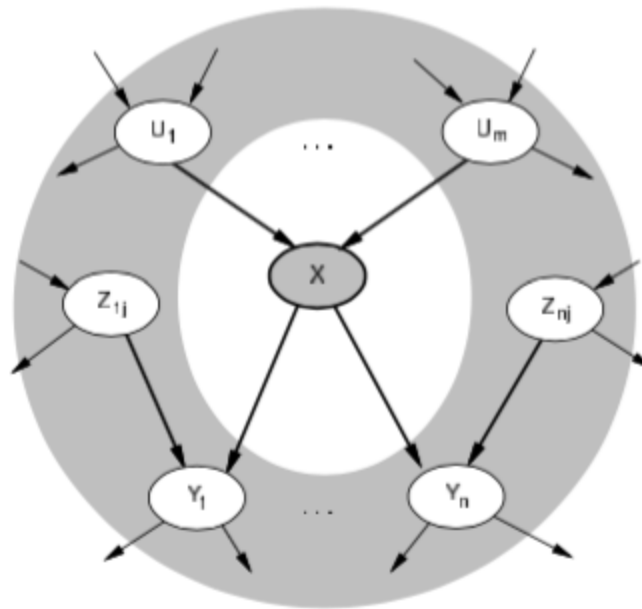
Local semantics



- Each node is conditionally independent of its nondescendants given its parents
- Theorem: Local semantics \Leftrightarrow global semantics

Markov blanket

- Each node is conditionally independent of all others given its **Markov blanket**
parents + children + children's parents



Constructing Bayesian networks

- A series of locally testable assertions of conditional independence
- Guarantees the required global semantics
- Algorithm
 - Choose an ordering of variables X_1, \dots, X_n
 - For $i = 1$ to n
 - add X_i to the network
 - select parents from X_1, \dots, X_{i-1} such that
$$\mathbf{P}(X_i | \text{Parents}(X_i)) = \mathbf{P}(X_i | X_1, \dots, X_{i-1})$$

Constructing Bayesian networks

- The choice of parents guarantees the global semantics:

$$\begin{aligned} \mathbf{P}(X_1, \dots, X_n) &= \prod_{i=1}^n \mathbf{P}(X_i | X_1, \dots, X_{i-1}) && \text{Chain rule} \\ &= \prod_{i=1}^n \mathbf{P}(X_i | \text{Parents}(X_i)) && \text{By construction} \end{aligned}$$

Example

- Suppose we choose the ordering M, J, A, B, E

$$P(J|M)=P(J)?$$

No

$$P(A|J, M)=P(A|J)?$$

No

$$P(A|J, M)=P(A)?$$

Yes

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

No

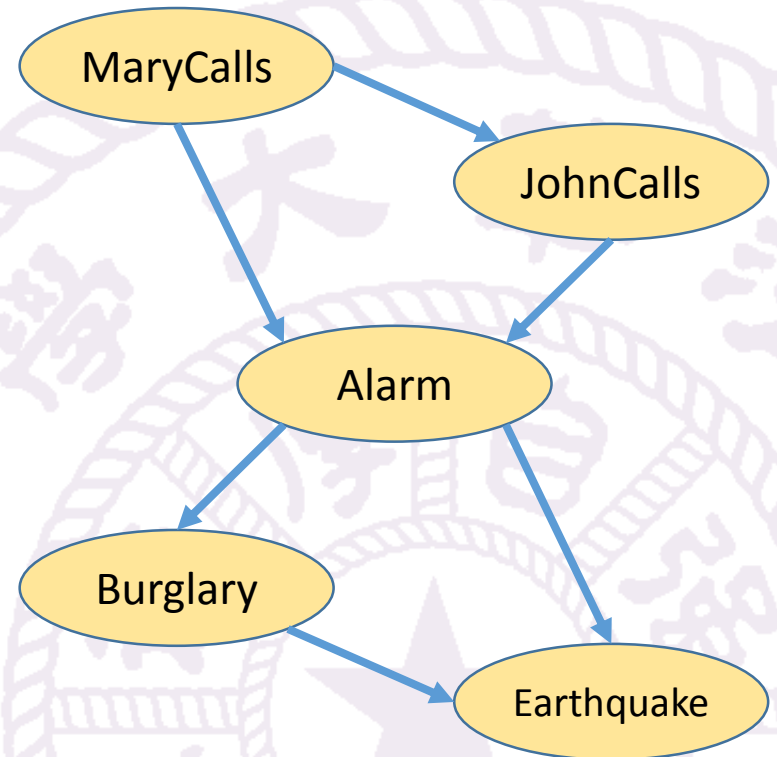
$$P(B|A, J, M)=P(B)?$$

Yes

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

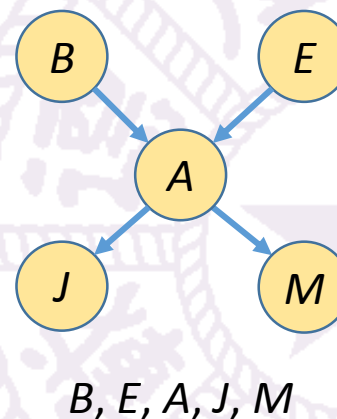
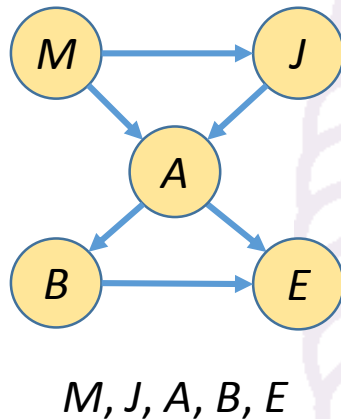
No

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

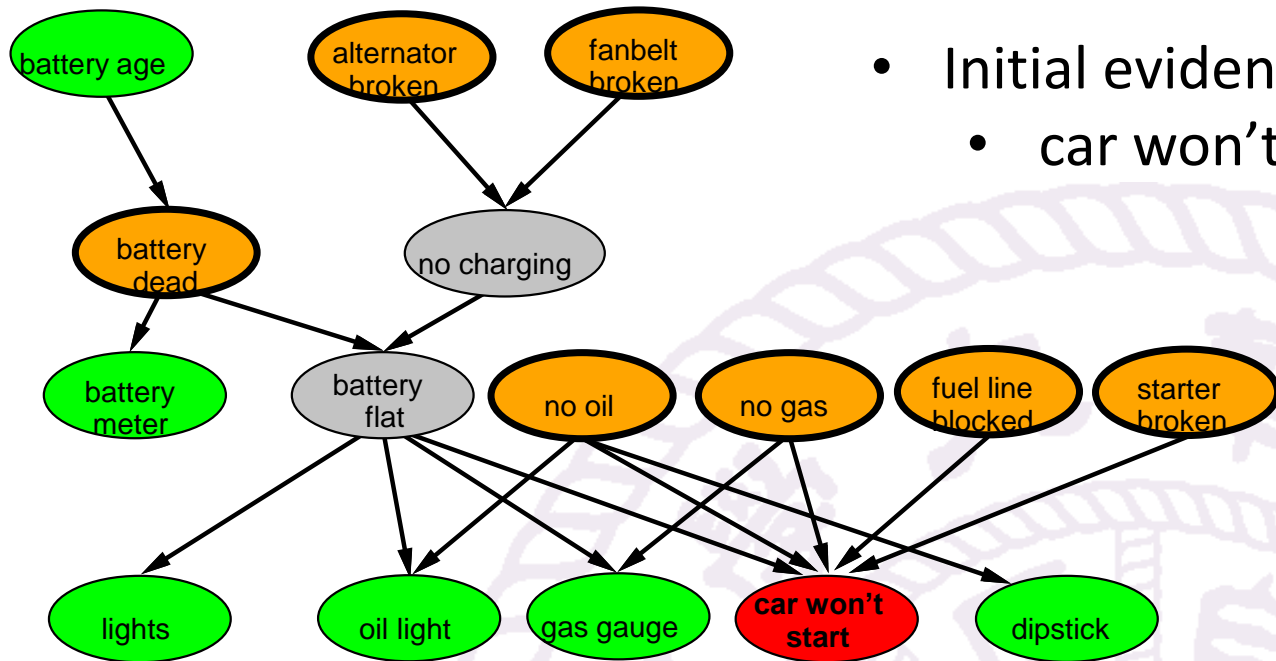


Example

- Deciding conditional independence is hard in noncausal directions
- Assessing conditional probabilities is hard in noncausal directions
- Network is less compact: $1+2+4+2+4=13$ numbers needed

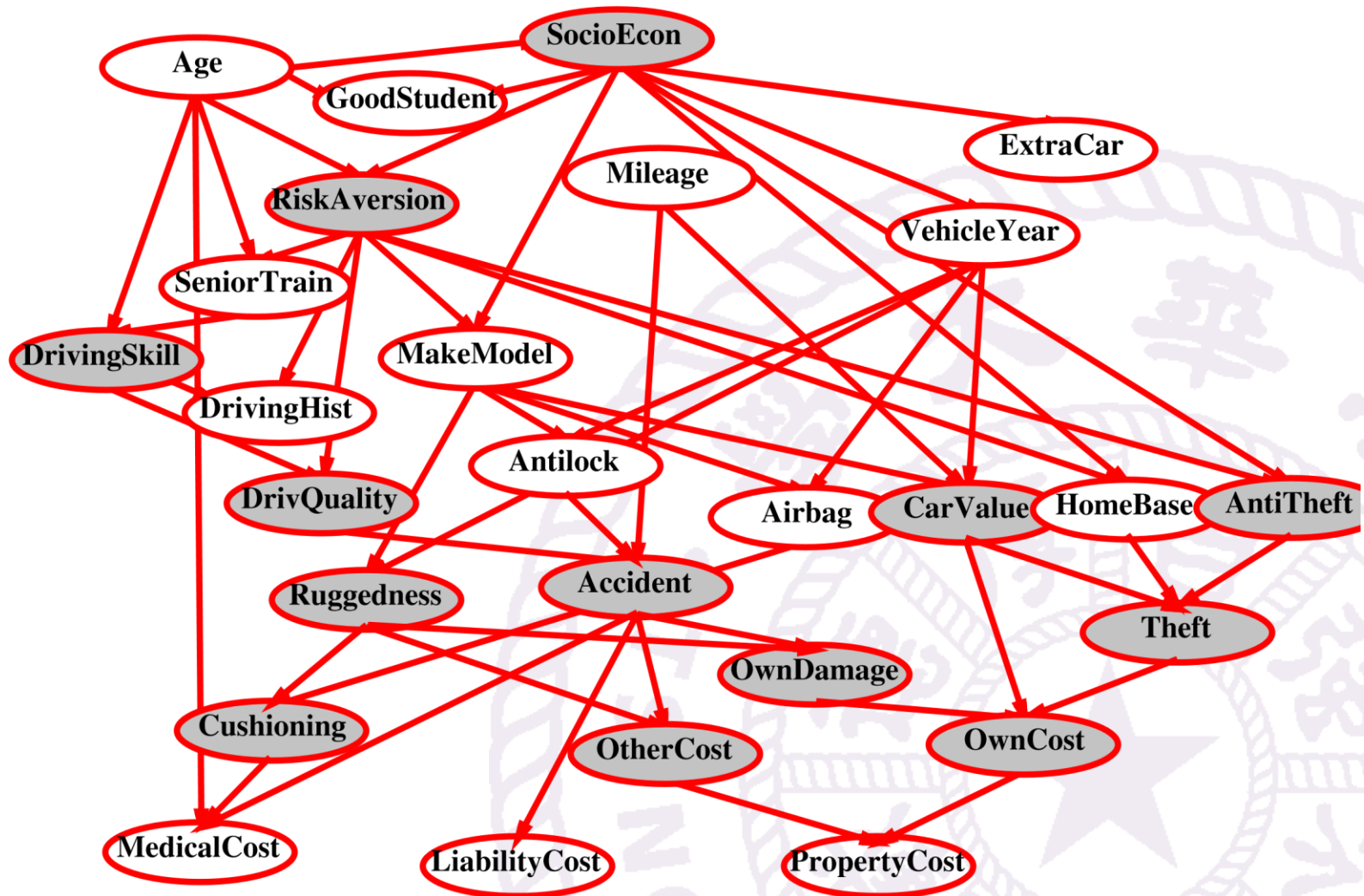


Example: Car diagnosis



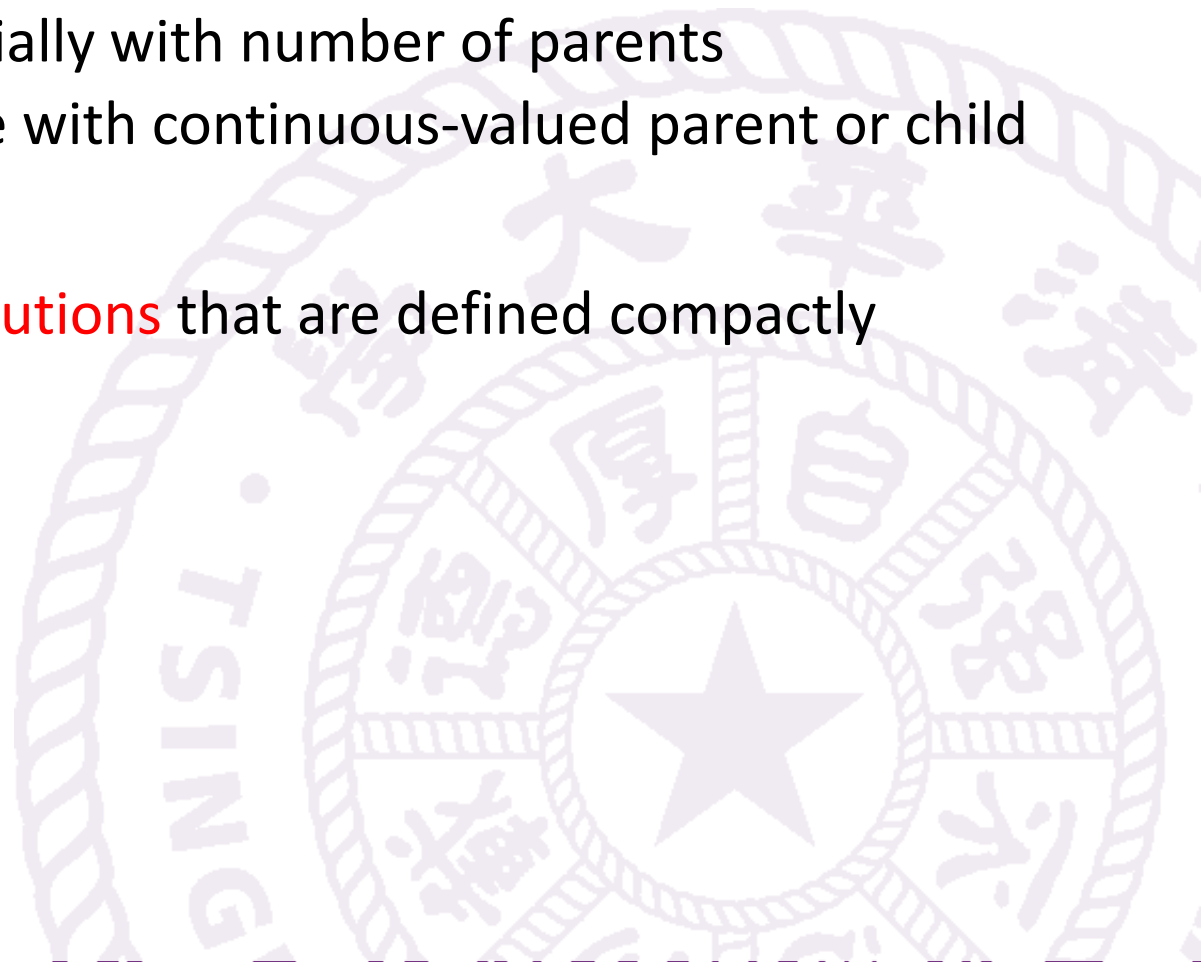
- Testable variables (green)
- “broken, so fix it” variables (orange)
- Hidden variables (gray)
 - ensure sparse structure, reduce parameters

Example: Car insurance



Compact conditional distributions

- CPT
 - grows exponentially with number of parents
 - becomes infinite with continuous-valued parent or child
- Solution
 - **canonical distributions** that are defined compactly



Deterministic models

- the simplest case

$$X = f(\text{Parents}(X)) \text{ for some function } f$$

- e.g., Boolean functions

$$\text{NorthAmerican} \Leftrightarrow \text{Canadian} \vee \text{US} \vee \text{Mexican}$$

- e.g., numerical relationships among continuous variables

$$\frac{\partial \text{Level}}{\partial t} = \text{inflow} + \text{precipitation} - \text{outflow} - \text{evaporation}$$

Noisy-OR distributions

- Multiple noninteracting causes
 - Parents U_1, \dots, U_k include all causes (can add **leak node**)
 - Negated $\neg U_i$ causes do not have any influence on X
 - Independent failure probability q_i for each cause alone

$$P(\neg X | U_1, \dots, U_j, \neg U_{j+1}, \dots, \neg U_k) = \prod_{i=1}^j q_i$$

$$P(X | U_1, \dots, U_j, \neg U_{j+1}, \dots, \neg U_k) = 1 - \prod_{i=1}^j q_i$$

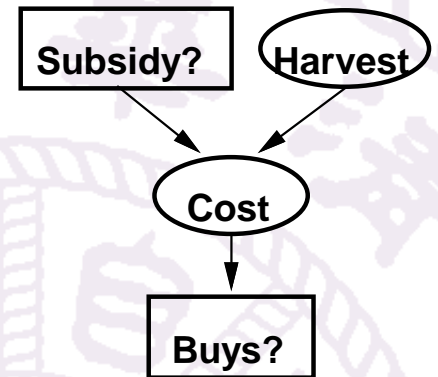
- Number of parameters **linear** in number of parents

Noisy-OR distributions

<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	$P(\text{Fever})$	$P(\neg \text{Fever})$
F	F	F	0.0	1.0
F	F	T	0.9	0.1
F	T	F	0.8	0.2
F	T	T	0.98	$0.02=0.2 \times 0.1$
T	F	F	0.4	0.6
T	F	T	0.94	$0.06=0.6 \times 0.1$
T	T	F	0.88	$0.12=0.6 \times 0.2$
T	T	T	0.988	$0.012=0.6 \times 0.2 \times 0.1$

Hybrid (discrete+continuous) networks

- Discrete (Subsidy? and Buys?)
- Continuous (Harvest and Cost)
- Option 1: discretization—possibly large errors, large CPTs
- Option 2: finitely parameterized canonical families
 - Continuous variable, discrete+continuous parents (e.g., Cost)
 - Discrete variable, continuous parents (e.g., Buys?)



Continuous child variables

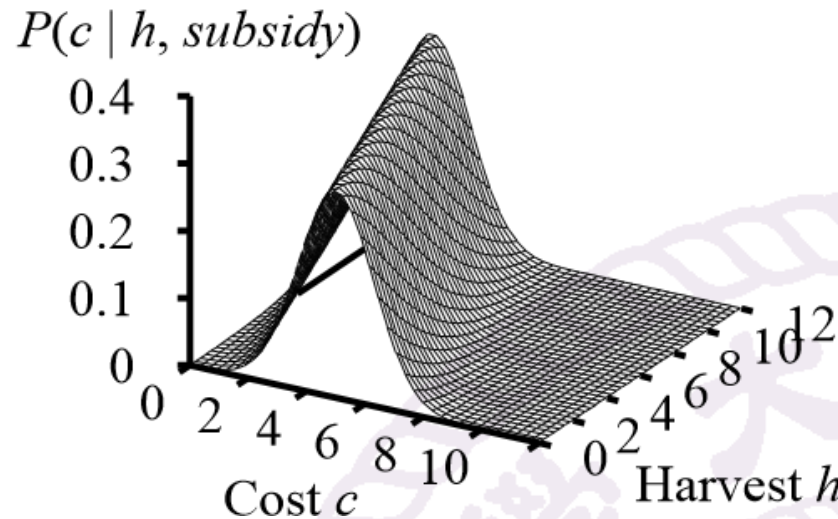
- One **conditional density function** for child variable given continuous parents, for each possible assignment to discrete parents

- **Linear Gaussian model**, e.g.,

$$\begin{aligned} & p(\text{Cost} = c | \text{Harvest} = h, \text{Subsidy} = \text{true}) \\ &= N(c | a_t h + b_t, \sigma_t) \\ &= \frac{1}{\sigma_t \sqrt{2\pi}} e^{\left(-\frac{1}{2} \left(\frac{c - (a_t h + b_t)}{\sigma_t} \right)^2 \right)} \end{aligned}$$

- Mean **Cost** varies linearly with **Harvest**, variance is fixed

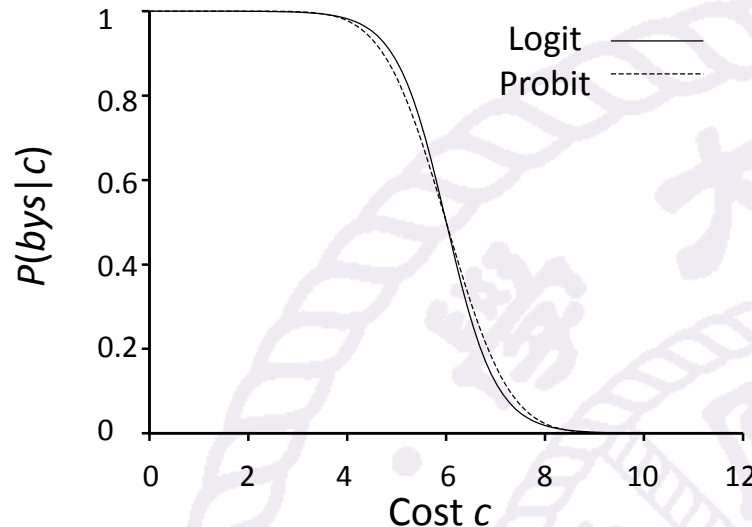
Continuous child variables



- All-continuous network with LG distributions
 - full joint distribution is a multivariate Gaussian
- Discrete+continuous LG network is a **conditional Gaussian network**
 - a multivariate Gaussian over all continuous variables for each combination of discrete variable values

Discrete variable w/ continuous parents

- Probability of *Buys?* given *Cost* should be a “soft” threshold



- **Probit** distribution uses integral of Gaussian:

$$\Phi(x) = \int_{-\infty}^x N(0,1)(x)dx$$

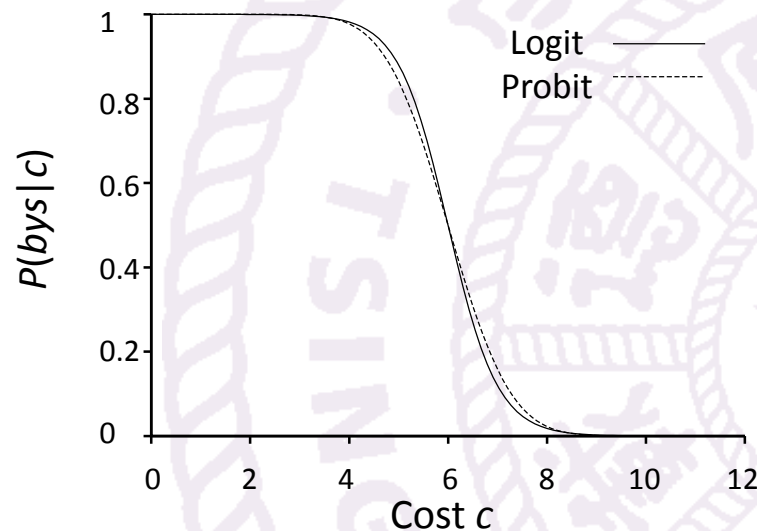
$$P(Buys = true|Cost = c) = \Phi((-c + \mu)/\sigma)$$

Discrete variable w/ continuous parents

- **Sigmoid** (or **logit**) distribution also used in neural networks

$$P(\text{Buys} = \text{true} | \text{Cost} = c) = \frac{1}{1 + e^{-2\frac{-c+\mu}{\sigma}}}$$

- Sigmoid has similar shape to probit but much longer tails:



Inference tasks

- Simple queries

- compute posterior marginal $P(X_i | E=e)$
- e.g., $P(\text{NoGas} | \text{Gauge}=\text{empty}, \text{Lights}=\text{on}, \text{Starts}=\text{false})$

- Conjunctive queries

$$P(X_i, X_j | E=e) = P(X_i | E=e)P(X_j | X_i, E=e)$$

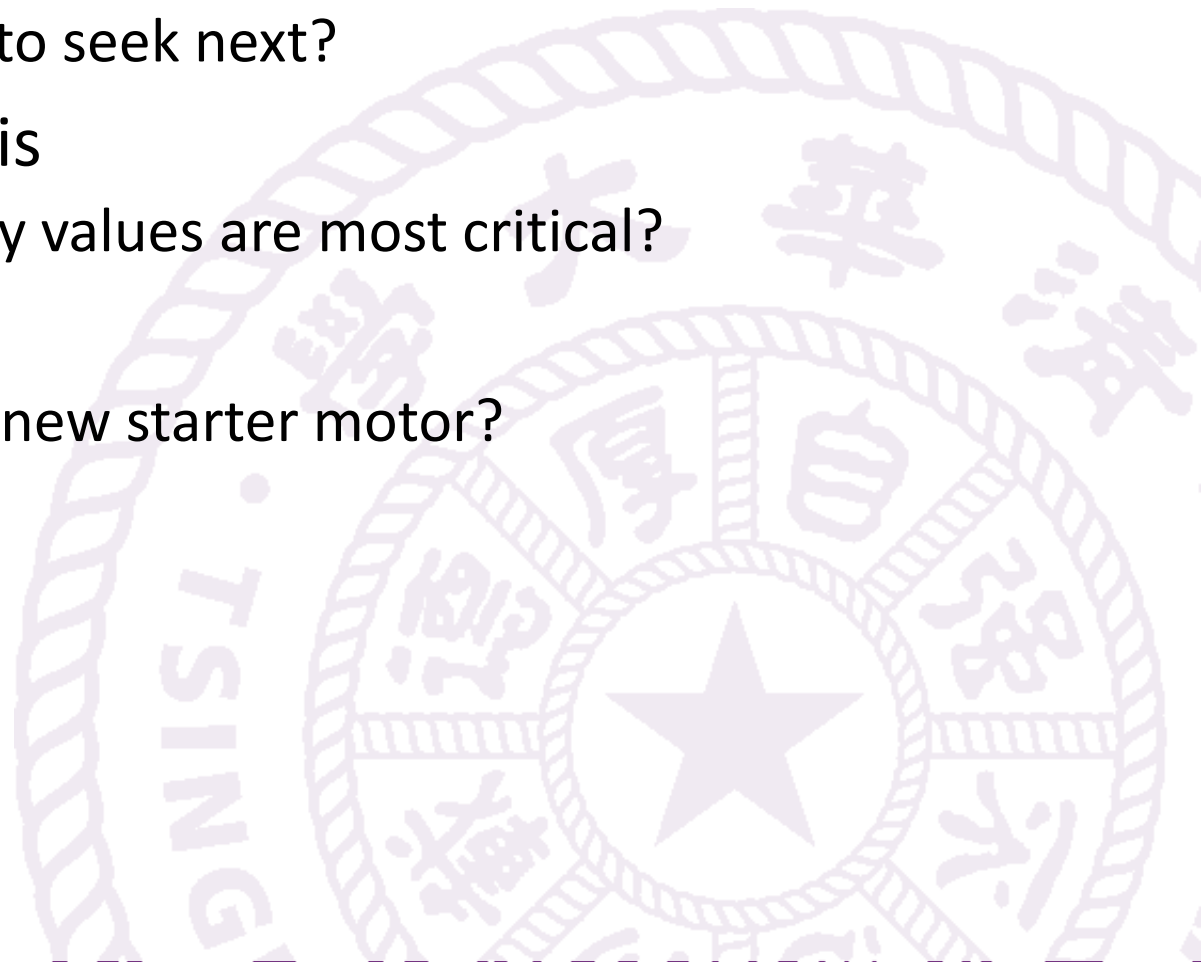
- Optimal decisions

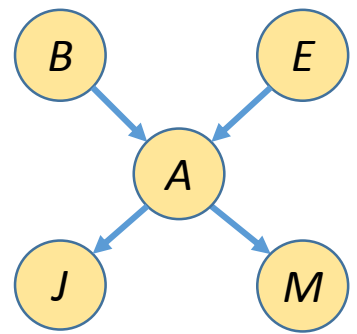
- decision networks include utility information
- probabilistic inference required for

$$P(\text{outcome} | \text{action}, \text{evidence})$$

Inference tasks

- Value of information
 - which evidence to seek next?
- Sensitivity analysis
 - which probability values are most critical?
- Explanation
 - why do I need a new starter motor?





Inference by enumeration

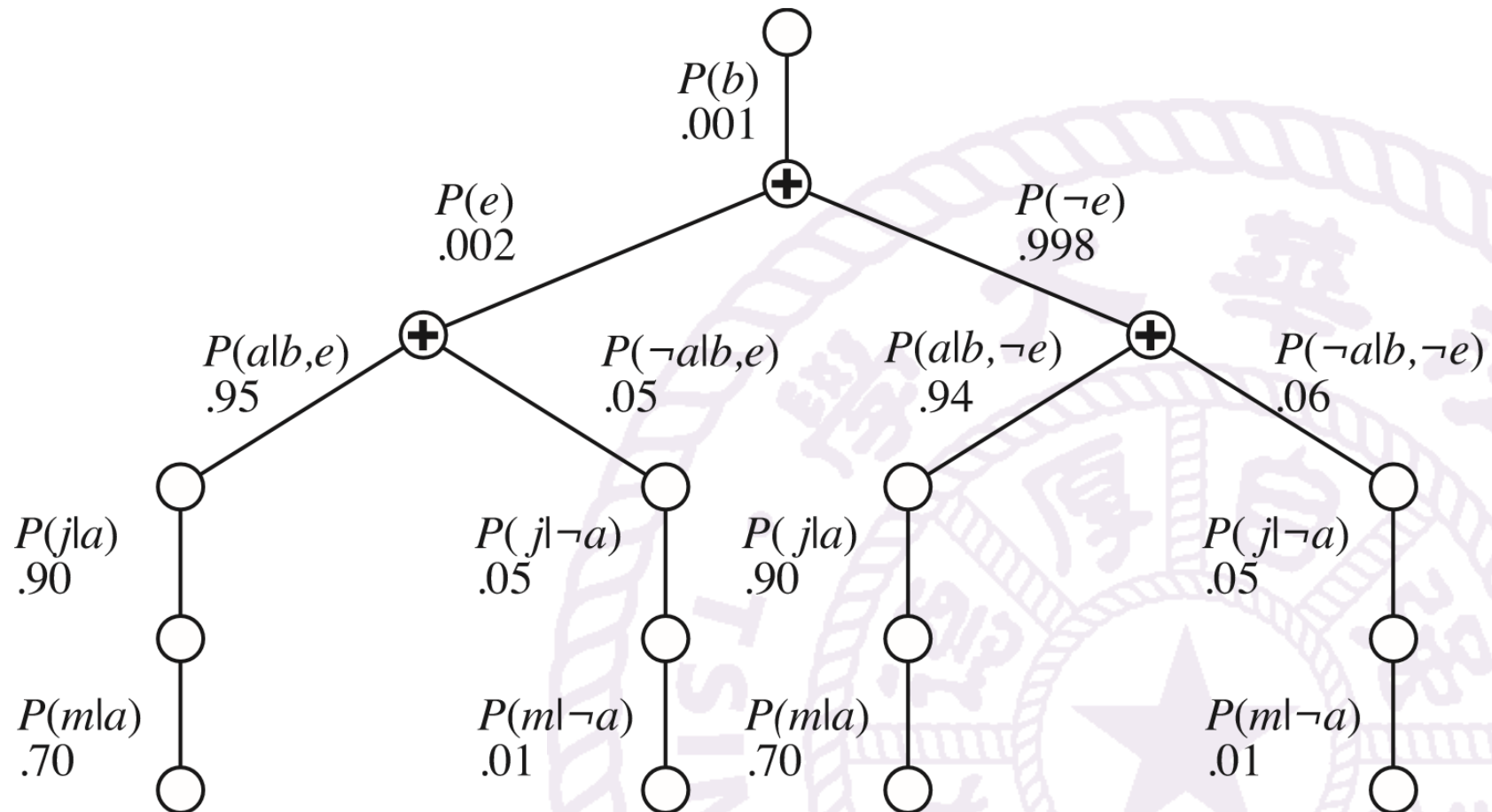
- sum out variables from the joint without actually constructing its explicit representation
- Simple query on the burglary network:

$$\begin{aligned}\mathbf{P}(B|j, m) &= \mathbf{P}(B, j, m) / \mathbf{P}(j, m) = \alpha \mathbf{P}(B, j, m) \\ &= \alpha \sum_e \sum_a \mathbf{P}(B, e, a, j, m)\end{aligned}$$

- Rewrite full joint entries using product of CPT entries

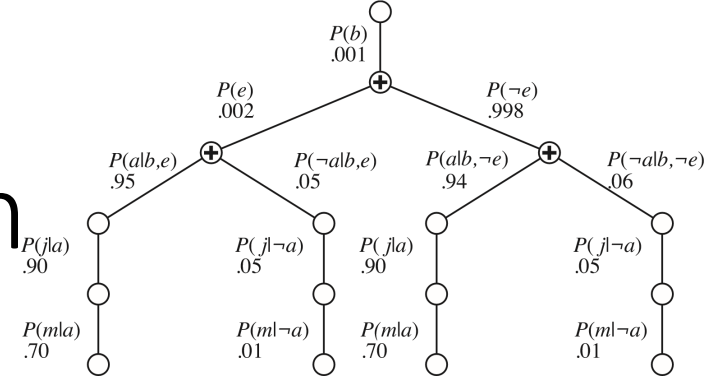
$$\begin{aligned}\mathbf{P}(B|j, m) &= \alpha \sum_e \sum_a \mathbf{P}(B)P(e)\mathbf{P}(a|B, e)P(j|a)P(m|a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e)P(j|a)P(m|a)\end{aligned}$$

Evaluation tree



$$P(b|j, m) = \alpha P(b) \sum_e P(e) \sum_a P(a|b, e) P(j|a) P(m|a)$$

Enumeration algorithm



function ENUMERATION-ASK(X, \mathbf{e}, bn) **returns** a distribution over X

inputs: X , the query variable

\mathbf{e} , observed values for variables \mathbf{E}

bn , a Bayes net with variables $\{X\} \cup \mathbf{E} \cup \mathbf{Y}$ /* $\mathbf{Y} = \text{hidden variables}$ */

$\mathbf{Q}(X) \leftarrow$ a distribution over X , initially empty

for each value x_i of X **do**

$\mathbf{Q}(x_i) \leftarrow$ ENUMERATE-ALL($bn.VARS, \mathbf{e}_{x_i}$)

where \mathbf{e}_{x_i} is \mathbf{e} extended with $X = x_i$

return NORMALIZE($\mathbf{Q}(X)$)

function ENUMERATE-ALL($vars, \mathbf{e}$) **returns** a real number

if EMPTY?($vars$) **then return** 1.0

$Y \leftarrow$ FIRST($vars$)

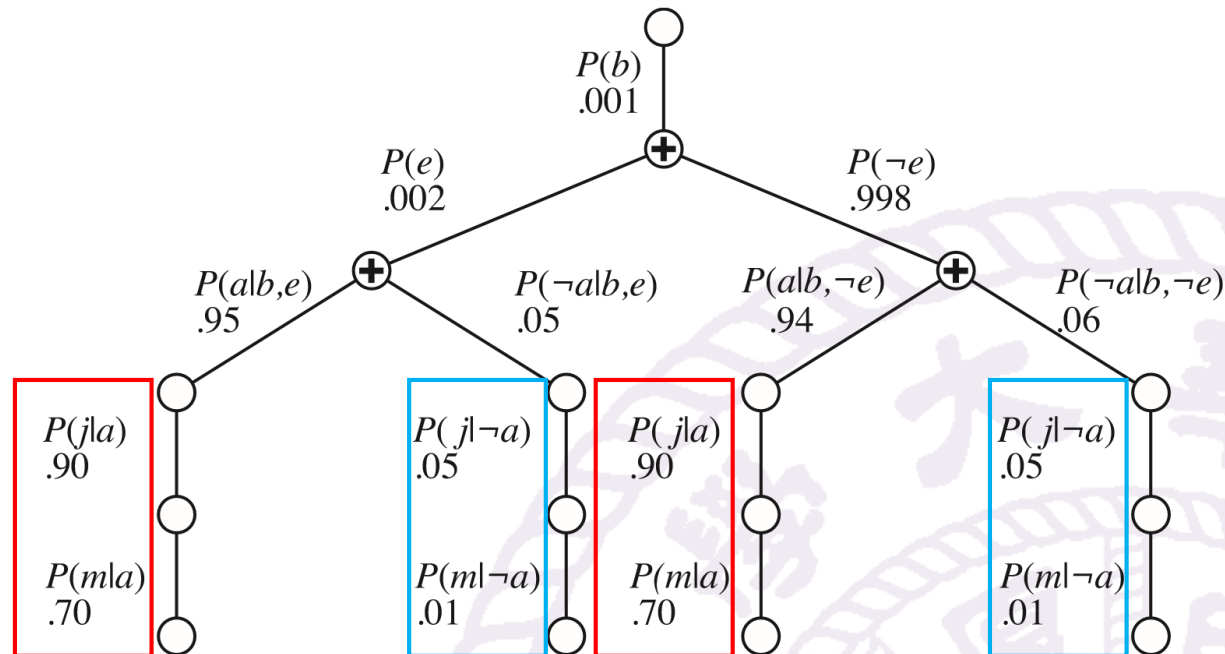
if Y has value y in \mathbf{e}

then return $P(y \mid \text{parents}(Y)) \times \text{ENUMERATE-ALL}(\text{REST}(vars), \mathbf{e})$

else return $\sum_y P(y \mid \text{parents}(Y)) \times \text{ENUMERATE-ALL}(\text{REST}(vars), \mathbf{e}_y)$

where \mathbf{e}_y is \mathbf{e} extended with $Y = y$

Inference by enumeration



- Recursive depth-first enumeration
 - $O(n)$ space, $O(d^n)$ time
- Enumeration is inefficient: repeated computation
 - e.g., computes $P(j|a)P(m|a)$ for each value of E

Inference by variable elimination

- Variable elimination: carry out summations right-to-left, storing intermediate results (**factors**) to avoid recomputation

$$P(B|j, m) = \underbrace{\alpha P(B)}_B \sum_e \underbrace{P(e)}_E \sum_a \underbrace{P(a|B, e)}_A \underbrace{P(j|a)}_J \underbrace{P(m|a)}_M$$

$$= \alpha P(B) \sum_e P(e) \sum_a P(a|B, e) P(j|a) f_M(a)$$

$$= \alpha P(B) \sum_e P(e) \sum_a P(a|B, e) f_J(a) f_M(a)$$

$$= \alpha P(B) \sum_e P(e) \sum_a f_A(a, b, e) f_J(a) f_M(a)$$

$$\text{sum}_{\text{out } A} = \alpha P(B) \sum_e f_E(e) f_{\bar{A}JM}(b, e)$$

$$\text{sum}_{\text{out } E} = \alpha P(B) f_{\bar{E}\bar{A}JM}(b)$$

$$= \alpha f_B(b) \times f_{\bar{E}\bar{A}JM}(b)$$

Variable elimination: Basic operations

- **Pointwise product** of factors f_1 and f_2 :

$$\begin{aligned} f_1(x_1, \dots, x_j, y_1, \dots, y_k) \times f_2(y_1, \dots, y_k, z_1, \dots, z_l) \\ = f(x_1, \dots, x_j, y_1, \dots, y_k, z_1, \dots, z_l) \end{aligned}$$

- E.g., $f_1(a, b) \times f_2(b, c) = f(a, b, c)$

Example

A	B	$f_1(A, B)$	B	C	$f_2(B, C)$	A	B	C	$f(A, B, C)$
T	T	0.3	T	T	0.2	T	T	T	$0.3 \cdot 0.2$
T	F	0.7	T	F	0.8	T	T	F	$0.3 \cdot 0.8$
F	T	0.9	F	T	0.6	T	F	T	$0.7 \cdot 0.6$
F	F	0.1	F	F	0.4	T	F	F	$0.7 \cdot 0.4$
						F	T	T	$0.9 \cdot 0.2$
						F	T	F	$0.9 \cdot 0.8$
						F	F	T	$0.1 \cdot 0.6$
						F	F	F	$0.1 \cdot 0.4$

Variable elimination: Basic operations

- **Summing out** a variable from a product of factors:
 - move any constant factors outside the summation
 - add up submatrices in pointwise product of remaining factors

$$\begin{aligned} & \sum_X f_1 \times \cdots \times f_k \\ &= f_1 \times \cdots \times f_i \sum_X f_{i+1} \times \cdots \times f_k \\ &= f_1 \times \cdots \times f_i \times f_{\bar{X}} \end{aligned}$$

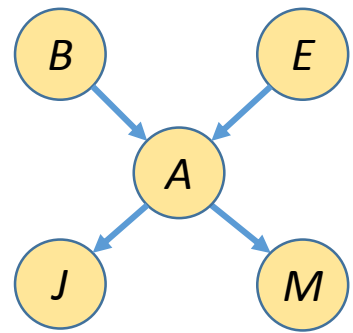
assuming f_1, \dots, f_i do not depend on X

Example

$$\begin{aligned} & \mathbf{f}(B, C) \\ = & \sum_a \mathbf{f}_3(A, B, C) \\ = & \mathbf{f}_3(a, B, C) + \mathbf{f}_3(\neg a, B, C) \\ = & \begin{pmatrix} 0.06 & 0.24 \\ 0.42 & 0.28 \end{pmatrix} + \begin{pmatrix} 0.18 & 0.72 \\ 0.06 & 0.04 \end{pmatrix} \\ = & \begin{pmatrix} 0.24 & 0.96 \\ 0.48 & 0.32 \end{pmatrix} \end{aligned}$$

Variable elimination algorithm

```
function ELIMINATION-ASK( $X, \mathbf{e}, bn$ ) returns a distribution over  $X$   
  inputs:  $X$ , the query variable  
            $\mathbf{e}$ , observed values for variables  $\mathbf{E}$   
            $bn$ , a Bayesian network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$   
  
   $factors \leftarrow []$   
  for each  $var$  in ORDER( $bn.VARS$ ) do  
     $factors \leftarrow [MAKE-FACTOR(var, \mathbf{e}) | factors]$   
    if  $var$  is a hidden variable then  $factors \leftarrow SUM-OUT(var, factors)$   
  return NORMALIZE(POINTWISE-PRODUCT( $factors$ ))
```



Irrelevant variables

- Consider the query $\mathbf{P(JohnCalls | Burglary = true)}$

$$\mathbf{P(J|b)}$$

$$= \alpha P(b) \sum_e P(e) \sum_a \mathbf{P(a|b, e)} P(J|a) \sum_m P(m|a)$$

- Sum over \mathbf{m} is identically 1; \mathbf{M} is **irrelevant** to the query

- Theorem:

Y is irrelevant unless $Y \in \text{Ancestors}(\{X\} \cup \mathbf{E})$

- $X = \text{JohnCalls}$, $\mathbf{E} = \{\text{Burglary}\}$, and $\text{Ancestors}(\{X\} \cup \mathbf{E}) = \{\text{Alarm}, \text{Earthquake}\}$, so MaryCalls is irrelevant

Complexity of exact inference

- Singly connected networks (or polytrees):
 - any two nodes are connected by at most one (undirected) path
 - Complexity of variable elimination are linear in CPT sizes
- Multiply connected networks:
 - can reduce 3-SAT to exact inference \Rightarrow NP-hard
 - equivalent to counting 3-SAT models \Rightarrow #P-complete

Summary

- Bayes nets provide a natural representation for (causally induced) conditional independence
- Topology + CPTs = compact representation of joint distribution
- Canonical distributions (e.g., noisy-OR) = compact representation of CPTs
- Continuous variables => parameterized distributions (e.g., linear Gaussian)
- Exact inference by variable elimination

谢谢！

