

Course Overview

Lecture 9 Bayesian Networks

Marco Chiarandini

Department of Mathematics & Computer Science
University of Southern Denmark

Slides by Stuart Russell and Peter Norvig

- ✓ Introduction
 - ✓ Artificial Intelligence
 - ✓ Intelligent Agents
- ✓ Search
 - ✓ Uninformed Search
 - ✓ Heuristic Search
- ✓ Adversarial Search
 - ✓ Minimax search
 - ✓ Alpha-beta pruning
- ✓ Knowledge representation and Reasoning
 - ✓ Propositional logic
 - ✓ First order logic
 - ✓ Inference
- Uncertain knowledge and Reasoning
 - Probability and Bayesian approach
 - Bayesian Networks
 - Hidden Markov Chains
 - Kalman Filters
- Learning
 - Decision Trees
 - Maximum Likelihood
 - EM Algorithm
 - Learning Bayesian Networks
 - Neural Networks
 - Support vector machines

2

Outline

Probability Basis
Bayesian networks
Inference in BN

1. Probability Basis
2. Bayesian networks
3. Inference in BN

Summary

Probability Basis
Bayesian networks
Inference in BN

- Interpretations of probability
- Axioms of Probability
- (Continuous/Discrete) Random Variables
- Prior probability, joint probability, conditional or posterior probability, chain rule
- Inference by enumeration

How to reduce the computation of inference?

DEFINITION

INDEPENDENT EVENTS

Two events A and B are *independent* of each other if and only if $p(A \cap B) = p(A) p(B)$. When $p(B) \neq 0$ this is the same as saying that $p(A) = p(A|B)$. That is, knowing that B is true does not affect the probability of A being true.

CONDITIONALLY INDEPENDENT EVENTS

Two events A and B are said to be *conditionally independent* of each other, given event C if and only if $p(A \cap B | C) = p(A | C) p(B | C)$.

Conditional independence

$P(\textit{Toothache}, \textit{Cavity}, \textit{Catch})$ has $2^3 - 1 = 7$ independent entries

If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:

$$(1) P(\textit{catch}|\textit{toothache}, \textit{cavity}) = P(\textit{catch}|\textit{cavity})$$

The same independence holds if I haven't got a cavity:

$$(2) P(\textit{catch}|\textit{toothache}, \neg\textit{cavity}) = P(\textit{catch}|\neg\textit{cavity})$$

Catch is **conditionally independent** of *Toothache* given *Cavity*:

$$P(\textit{Catch}|\textit{Toothache}, \textit{Cavity}) = P(\textit{Catch}|\textit{Cavity})$$

Equivalent statements:

$$P(\textit{Toothache}|\textit{Catch}, \textit{Cavity}) = P(\textit{Toothache}|\textit{Cavity})$$

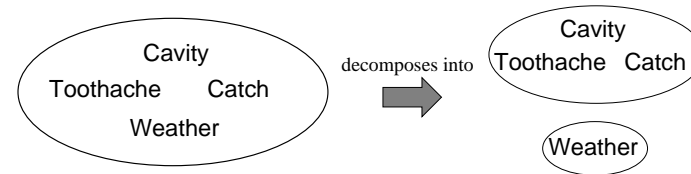
$$P(\textit{Toothache}, \textit{Catch}|\textit{Cavity}) = P(\textit{Toothache}|\textit{Cavity})P(\textit{Catch}|\textit{Cavity})$$

5

Independence

A and B are **independent** iff

$$P(A|B) = P(A) \quad \text{or} \quad P(B|A) = P(B) \quad \text{or} \quad P(A, B) = P(A)P(B)$$



$$P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}, \textit{Weather}) \\ = P(\textit{Toothache}, \textit{Catch}, \textit{Cavity})P(\textit{Weather})$$

32 entries reduced to 12; for n independent biased coins, $2^n \rightarrow n$

Absolute independence powerful but rare

Dentistry is a large field with hundreds of variables, none of which are independent. What to do?

6

Conditional independence contd.

Write out full joint distribution using chain rule:

$$P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}) \\ = P(\textit{Toothache}|\textit{Catch}, \textit{Cavity})P(\textit{Catch}, \textit{Cavity}) \\ = P(\textit{Toothache}|\textit{Catch}, \textit{Cavity})P(\textit{Catch}|\textit{Cavity})P(\textit{Cavity}) \\ = P(\textit{Toothache}|\textit{Cavity})P(\textit{Catch}|\textit{Cavity})P(\textit{Cavity})$$

I.e., $2 + 2 + 1 = 5$ independent numbers (equations 1 and 2 remove 2)

In most cases, the use of conditional independence reduces the size of the representation of the joint distribution from exponential in n to linear in n .

Conditional independence is our most basic and robust form of knowledge about uncertain environments.

7

8

Bayes' Rule

Product rule $P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$

$$\implies \text{Bayes' rule } P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

or in distribution form

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \alpha P(X|Y)P(Y)$$

Useful for assessing **diagnostic** probability from **causal** probability:

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$

E.g., let M be meningitis, S be stiff neck:

$$P(m|s) = \frac{P(s|m)P(m)}{P(s)} = \frac{0.8 \times 0.0001}{0.1} = 0.0008$$

Note: posterior probability of meningitis still very small!

Summary

Probability is a rigorous formalism for uncertain knowledge
Joint probability distribution specifies probability of every **atomic event**
Queries can be answered by summing over atomic events
For nontrivial domains, we must find a way to reduce the joint size
Independence and **conditional independence** provide the tools

Bayes' Rule and conditional independence

$$\begin{aligned} &P(\text{Cavity}|\text{toothache} \wedge \text{catch}) \\ &= \alpha P(\text{toothache} \wedge \text{catch}|\text{Cavity})P(\text{Cavity}) \\ &= \alpha P(\text{toothache}|\text{Cavity})P(\text{catch}|\text{Cavity})P(\text{Cavity}) \end{aligned}$$

This is an example of a **naive Bayes** model:

$$P(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = P(\text{Cause}) \prod_i P(\text{Effect}_i|\text{Cause})$$



Total number of parameters is **linear** in n

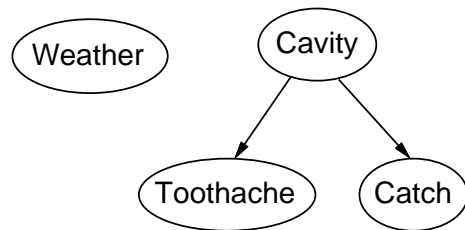
Outline

1. Probability Basis
2. Bayesian networks
3. Inference in BN

- ◇ Syntax
- ◇ Semantics
- ◇ Parameterized distributions

Example

Topology of network encodes conditional independence assertions:



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

13

15

Bayesian networks

Definition

A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions

Syntax:

- a set of nodes, one per variable
- a directed, acyclic graph (link \approx “directly influences”)
- a conditional distribution for each node given its parents:
 $P(X_i | Parents(X_i))$

In the simplest case, conditional distribution represented as a **conditional probability table** (CPT) giving the distribution over X_i for each combination of parent values

14

16

Example

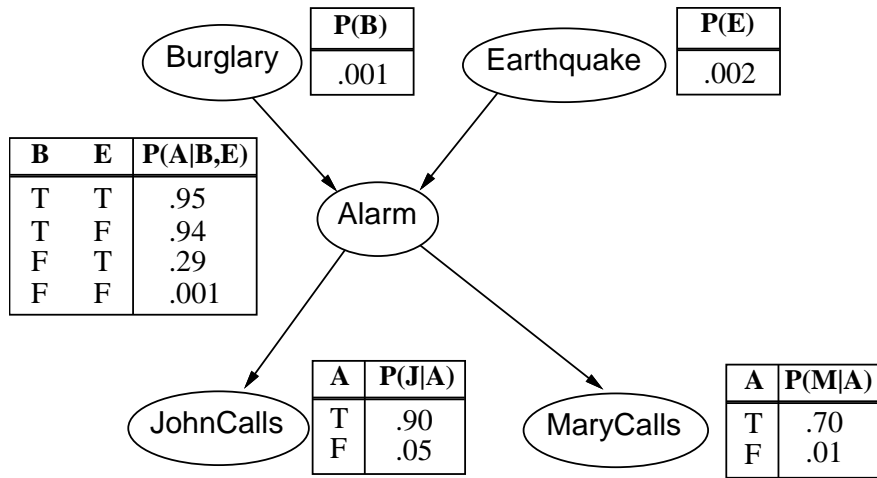
I’m at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn’t call. Sometimes it’s set off by minor earthquakes. Is there a burglar?

Variables: *Burglar*, *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 contd.



17

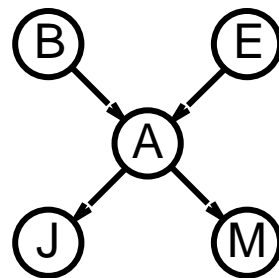
Global semantics

“Global” semantics defines the full joint distribution as the product of the local conditional distributions:

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

e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$\begin{aligned}
 &= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) \\
 &= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\
 &\approx 0.00063
 \end{aligned}$$



19

Compactness

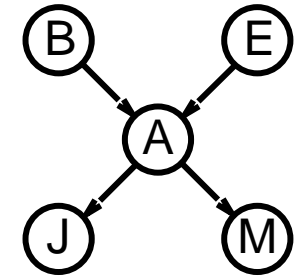
A CPT for Boolean X_i with k Boolean parents has 2^k rows for the combinations of parent values

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

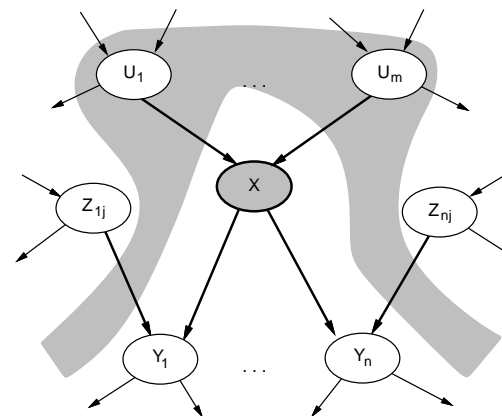
For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers (vs. $2^5 - 1 = 31$)



18

Local semantics

Local semantics: each node is conditionally independent of its nondescendants given its parents

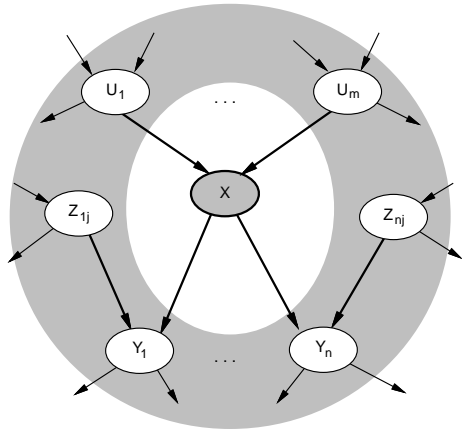


Theorem: Local semantics \Leftrightarrow global semantics

20

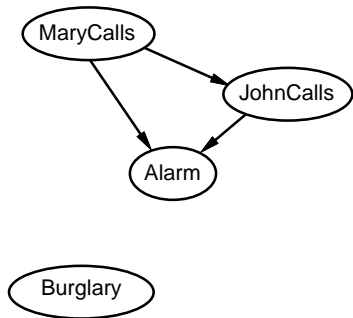
Markov blanket

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



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
 $P(A|J, M) = P(A)$? No
 $P(B|A, J, M) = P(B|A)$? Yes
 $P(B|A, J, M) = P(B)$? No
 $P(E|B, A, M) = P(E|A)$? No
 $P(E|B, A, J, M) = P(E|A, B)$? Yes
 Deciding conditional independence is hard in noncausal directions (Causal models and conditional independence seem hardwired for humans!)
 Assessing conditional probabilities is hard in noncausal directions
 Network is less compact:
 $1 + 2 + 4 + 2 + 4 = 13$ numbers needed

Constructing Bayesian networks

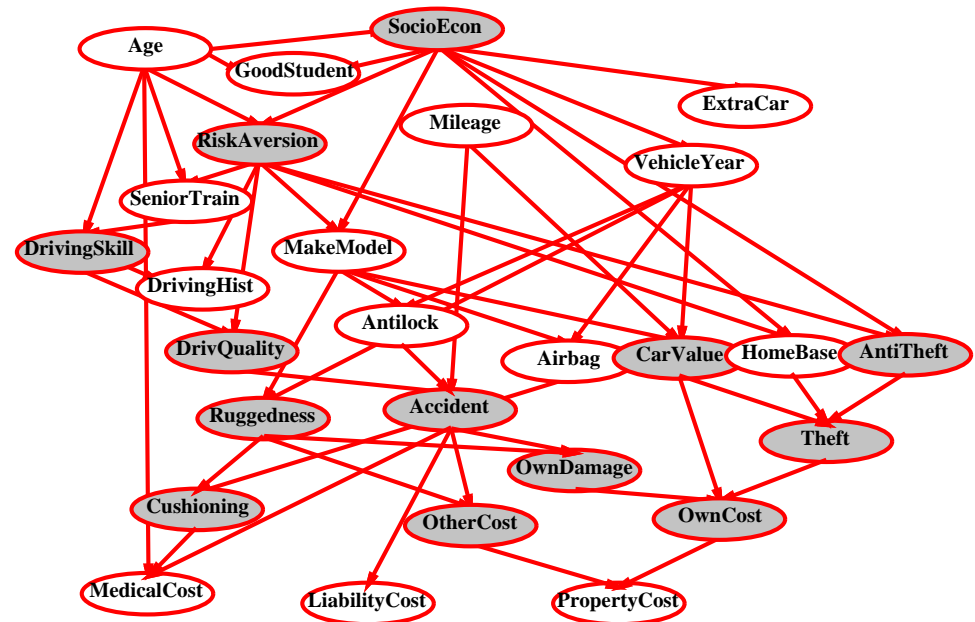
Need a method such that a series of locally testable assertions of conditional independence guarantees the required global semantics

- 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
 $P(X_i | Parents(X_i)) = P(X_i | X_1, \dots, X_{i-1})$

This choice of parents guarantees the global semantics:

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

Example: Car insurance



Compact conditional distributions

CPT grows exponentially with number of parents
CPT becomes infinite with continuous-valued parent or child

Solution:

canonical distributions that are defined compactly

Deterministic nodes are the simplest case:

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

E.g., Boolean functions

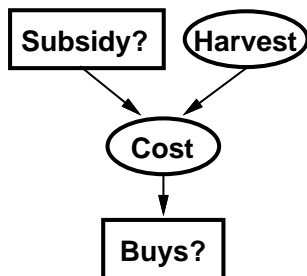
$$\text{NorthAmerican} \Leftrightarrow \text{Canadian} \vee \text{US} \vee \text{Mexican}$$

E.g., numerical relationships among continuous variables

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

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

- 1) Continuous variable, discrete+continuous parents (e.g., *Cost*)
- 2) Discrete variable, continuous parents (e.g., *Buys?*)

Compact conditional distributions contd.

Noisy-OR distributions model multiple noninteracting causes

- 1) Parents $U_1 \dots U_k$ include all causes (can add leak node)
- 2) Independent failure probability q_i for each cause alone

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

Cold	Flu	Malaria	P(Fever)	P(¬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 × 0.1
T	F	F	0.4	0.6
T	F	T	0.94	0.06 = 0.6 × 0.1
T	T	F	0.88	0.12 = 0.6 × 0.2
T	T	T	0.988	0.012 = 0.6 × 0.2 × 0.1

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

Continuous child variables

Need one conditional density function for child variable given continuous parents, for each possible assignment to discrete parents

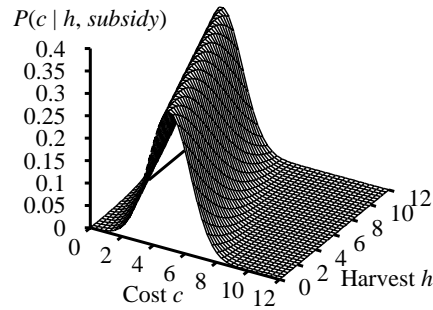
Most common is the linear Gaussian model, e.g.,:

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

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

↪ Linear variation is unreasonable over the full range but works OK if the **likely** range of *Harvest* is narrow

Continuous child variables

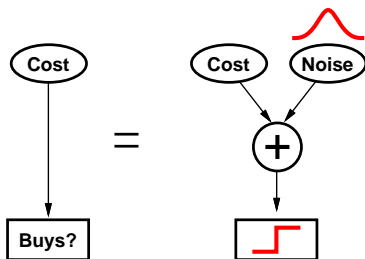


All-continuous network with linear Gaussian distributions
 \implies full joint distribution is a multivariate Gaussian

Discrete+continuous linear Gaussian network is a conditional Gaussian network i.e., a multivariate Gaussian over all continuous variables for each combination of discrete variable values

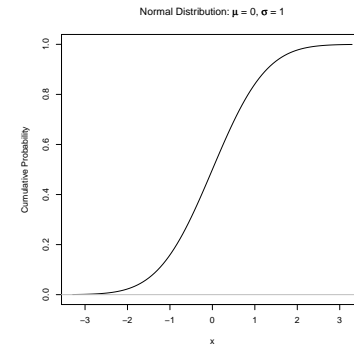
Why the probit?

1. It's sort of the right shape
2. Can be viewed as hard threshold whose location is subject to noise



Discrete variable w/ continuous parents

Probability of *Buys?* given *Cost* should be a "soft" threshold:



Probit distribution uses integral of Gaussian:

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

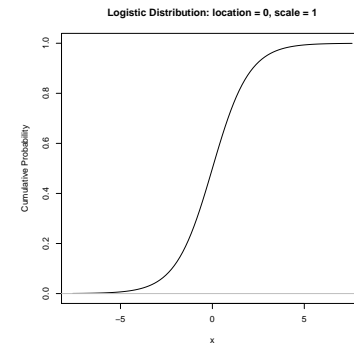
$$P(\text{Buys?} = \text{true} | \text{Cost} = c) = \Phi((-c + \mu)/\sigma)$$

Discrete variable contd.

Sigmoid (or logit) distribution also used in neural networks:

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

Sigmoid has similar shape to probit but much longer tails:



- Bayes nets provide a natural representation for (causally induced) conditional independence
- Topology + CPTs = compact representation of joint distribution
- Generally easy for (non)experts to construct
- Canonical distributions (e.g., noisy-OR) = compact representation of CPTs
- Continuous variables \implies parameterized distributions (e.g., linear Gaussian)

34

Inference tasks

- **Simple queries:** compute posterior marginal $\mathbf{P}(X_i|\mathbf{E}=\mathbf{e})$
e.g., $P(\text{NoGas}|\text{Gauge} = \text{empty}, \text{Lights} = \text{on}, \text{Starts} = \text{false})$
- **Conjunctive queries:** $\mathbf{P}(X_i, X_j|\mathbf{E}=\mathbf{e}) = \mathbf{P}(X_i|\mathbf{E}=\mathbf{e})\mathbf{P}(X_j|X_i, \mathbf{E}=\mathbf{e})$
- **Optimal decisions:** decision networks include utility information; probabilistic inference required for $P(\text{outcome}|\text{action}, \text{evidence})$
- **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?

36

Outline

1. Probability Basis
2. Bayesian networks
3. Inference in BN

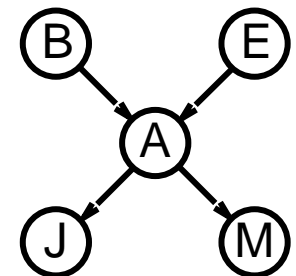
35

Inference by enumeration

Sum out variables from the joint without actually constructing its explicit representation

Simple query on the burglary network:

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



Rewrite full joint entries using product of CPT entries:

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

Recursive depth-first enumeration: $O(n)$ space, $O(d^n)$ time

37

Enumeration algorithm

```

function Enumeration-Ask( $X, e, bn$ ) returns a distribution over  $X$ 
  inputs:  $X$ , the query variable
             $e$ , observed values for variables  $E$ 
             $bn$ , a Bayesian network with variables  $\{X\} \cup E \cup Y$ 

   $Q(X) \leftarrow$  a distribution over  $X$ , initially empty
  for each value  $x_i$  of  $X$  do
    extend  $e$  with value  $x_i$  for  $X$ 
     $Q(x_i) \leftarrow$  Enumerate-All( $\text{Vars}[bn], e$ )
  return Normalize( $Q(X)$ )

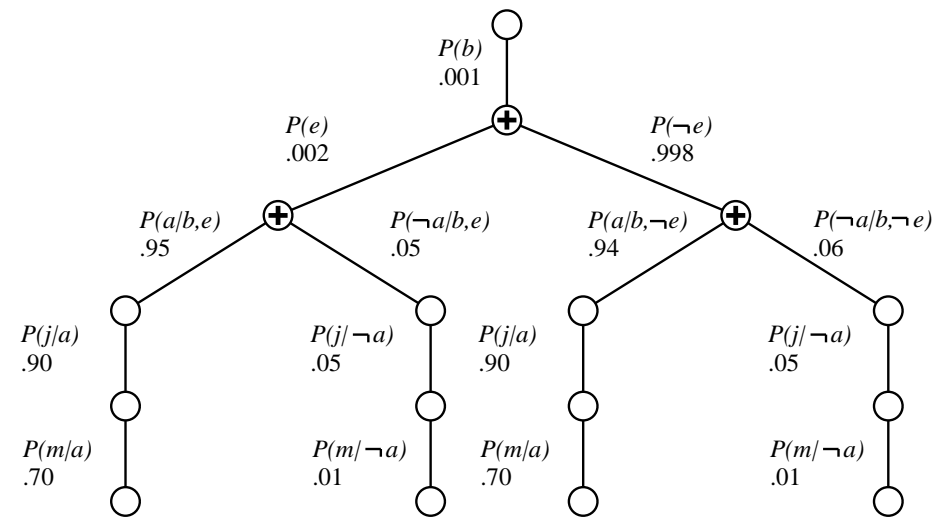


---


function Enumerate-All( $vars, e$ ) returns a real number
  if Empty?( $vars$ ) then return 1.0
   $Y \leftarrow$  First( $vars$ )
  if  $Y$  has value  $y$  in  $e$ 
    then return  $P(y \mid \text{parent}(Y)) \times$  Enumerate-All( $\text{Rest}(vars), e$ )
    else return  $\sum_y P(y \mid \text{parent}(Y)) \times$  Enumerate-All( $\text{Rest}(vars), e_y$ )
    where  $e_y$  is  $e$  extended with  $Y = y$ 
  
```

38

Evaluation tree



Enumeration is inefficient: repeated computation
e.g., computes $P(j|a)P(m|a)$ for each value of e

39

Complexity of exact inference

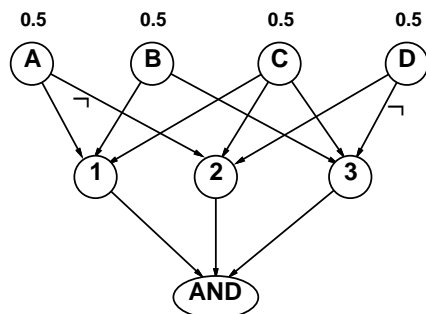
Singly connected networks (or polytrees):

- any two nodes are connected by at most one (undirected) path
- time and space cost (with variable elimination) are $O(d^k n)$
- hence time and space cost are linear in n and k bounded by a constant

Multiply connected networks:

- can reduce 3SAT to exact inference \implies NP-hard
- equivalent to counting 3SAT models \implies #P-complete

1. $A \vee B \vee C$
2. $C \vee D \vee \neg A$
3. $B \vee C \vee \neg D$



45

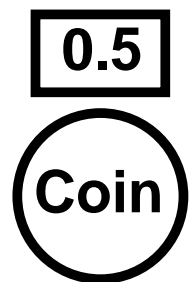
Inference by stochastic simulation

Basic idea:

- Draw N samples from a sampling distribution S
- Compute an approximate posterior probability \hat{P}
- Show this converges to the true probability P

Outline:

- Sampling from an empty network
- Rejection sampling: reject samples disagreeing with evidence
- Likelihood weighting: use evidence to weight samples
- Markov chain Monte Carlo (MCMC): sample from a stochastic process whose stationary distribution is the true posterior

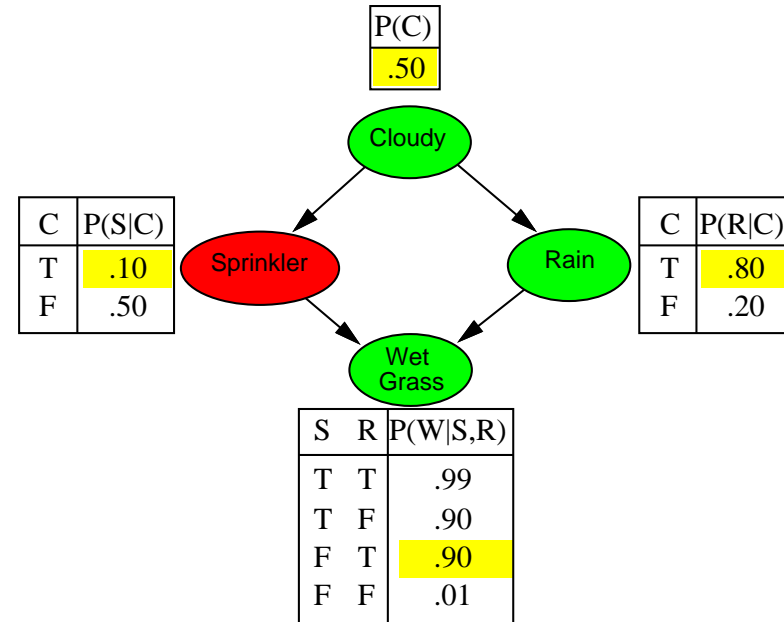


46

```

function Prior-Sample(bn) returns an event sampled from bn
  inputs: bn, a belief network specifying joint distribution
   $P(X_1, \dots, X_n)$ 
   $x \leftarrow$  an event with  $n$  elements
  for  $i = 1$  to  $n$  do
     $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
    given the values of  $\text{Parents}(X_i)$  in  $x$ 
  return  $x$ 
    
```

47



48

Sampling from an empty network contd.

Probability that `PriorSample` generates a particular event

$$S_{PS}(x_1 \dots x_n) = P(x_1 \dots x_n)$$

i.e., the true prior probability

E.g., $S_{PS}(t, f, t, t) = 0.5 \times 0.9 \times 0.8 \times 0.9 = 0.324 = P(t, f, t, t)$

Proof: Let $N_{PS}(x_1 \dots x_n)$ be the number of samples generated for event x_1, \dots, x_n . Then we have

$$\begin{aligned}
 \lim_{N \rightarrow \infty} \hat{P}(x_1, \dots, x_n) &= \lim_{N \rightarrow \infty} N_{PS}(x_1, \dots, x_n) / N \\
 &= S_{PS}(x_1, \dots, x_n) \\
 &= \prod_{i=1}^n P(x_i \mid \text{parents}(X_i)) = P(x_1 \dots x_n)
 \end{aligned}$$

⇒ That is, estimates derived from `PriorSample` are **consistent**

Shorthand: $\hat{P}(x_1, \dots, x_n) \approx P(x_1 \dots x_n)$

49