

Lecture of 11/16

Two parts:

Part 1: Practical issues in constructing Bayes networks

Part 2. Inference in Bayes networks

Part 1. Issues in constructing Bayes Nets

Constructing Belief Networks: Summary

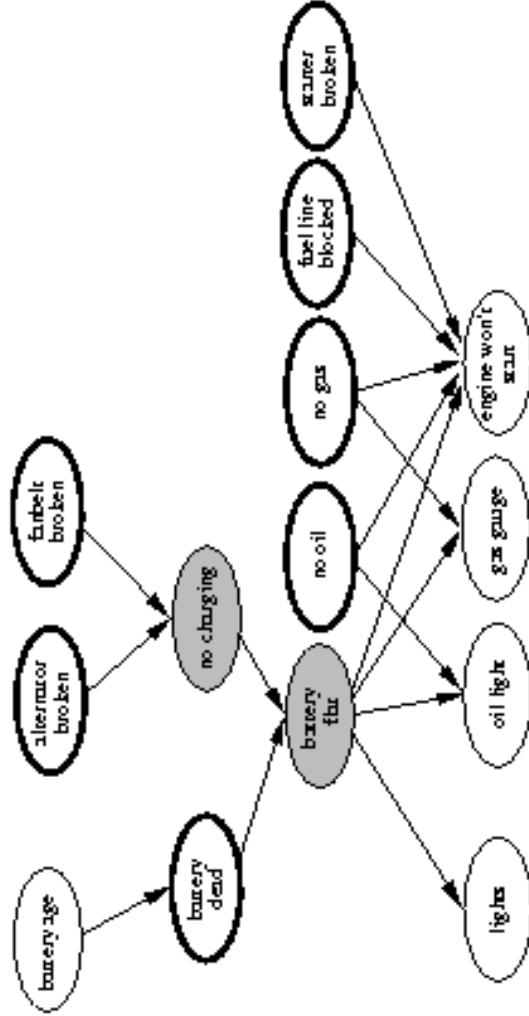
- [[Decide on what sorts of queries you are interested in answering
 - This in turn dictates what factors to model in the network
- Decide on a vocabulary of the variables and their domains for the problem
 - Introduce “Hidden” variables into the network as needed to make the network “sparse”
- Decide on an order of introduction of variables into the network
 - Introducing variables in causal direction leads to fewer connections (sparse structure) **AND** easier to assess probabilities
- Try to use canonical distributions to specify the CPTs
 - Noisy-OR
 - Parameterized discrete/continuous distributions
 - Such as Poisson, Normal (Gaussian) etc

Example: Car diagnosis

Initial evidence: engine won't start

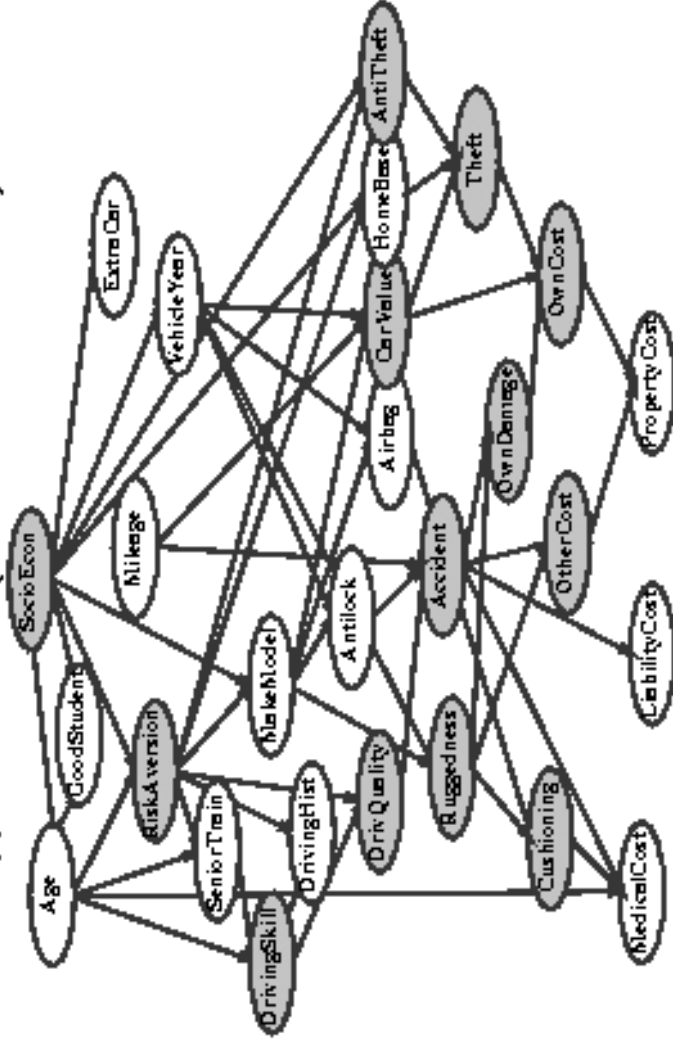
Testable variables (thin ovals), diagnosis variables (thick ovals)

Hidden variables (shaded) ensure sparse structure, reduce parameters



Example: Car insurance

Predict claim costs (medical, liability, property) given data on application form (other unshaded nodes)



Compact conditional distributions

CPT grows exponentially with no. 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

$$\text{North.American} \Leftrightarrow \text{Canadian} \vee \text{US} \vee \text{Mexican}$$

E.g., numerical relationships among continuous variables

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

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

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

<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	$0.012 = 0.6 \times 0.2 \times 0.1$

Number of parameters linear in number of parents

Constructing Belief Networks: Summary

- [[Decide on what sorts of queries you are interested in answering
 - This in turn dictates what factors to model in the network
- Decide on a vocabulary of the variables and their domains for the problem
 - Introduce “Hidden” variables into the network as needed to make the network “sparse”
- Decide on an order of introduction of variables into the network
 - Introducing variables in causal direction leads to fewer connections (sparse structure) **AND** easier to assess probabilities
- Try to use canonical distributions to specify the CPTs
 - Noisy-OR
 - Parameterized discrete/continuous distributions
 - Such as Poisson, Normal (Gaussian) etc

Case Study: Pathfinder System

- Domain: Lymph node diseases
 - Deals with 60 diseases and 100 disease findings
- Versions:
 - Pathfinder I: A rule-based system with logical reasoning
 - Pathfinder II: Tried a variety of approaches for uncertainty
 - Simple bayes reasoning outperformed
 - Pathfinder III: Simple bayes reasoning, but reassessed probabilities
 - Pathfinder IV: Bayesian network was used to handle a variety of conditional dependencies.
 - Deciding vocabulary: 8 hours
 - Devising the topology of the network: 35 hours
 - Assessing the (14,000) probabilities: 40 hours
 - Physician experts liked assessing causal probabilities
- Evaluation: 53 “referral” cases
 - Pathfinder III: 7.9/10
 - Pathfinder IV: 8.9/10 [Saves one additional life in every 1000 cases!]
 - A more recent comparison shows that Pathfinder now *outperforms* experts who helped design it!!

Part II. Inference in Bayes Nets

Inference tasks

Simple queries: compute posterior marginal $P(X_i | E = e)$

e.g., $P(NoGas | Gauge = empty, Lights = on, Start = false)$

Conjunctive queries: $P(X_i, X_j | E = e) = P(X_i | X_j, E = e)P(X_j | X_i, E = e)$

Optimal decisions: decision networks include utility information;

probabilistic inference required for $P(outcome | action, evidence)$

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?

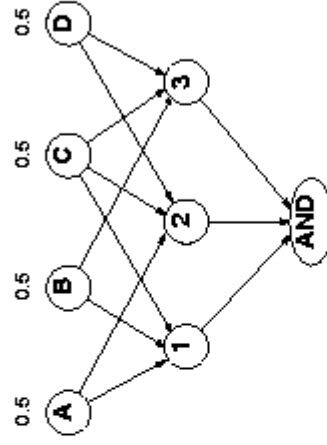
Complexity of exact inference

Singly connected networks (or polytrees):

- any two nodes are connected by at most one (undirected) path
- time and space cost of variable elimination are $O(d^{*n})$

Multiply connected networks:

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



1. A V B V C
2. C V D V \neg A
3. B V C V \neg D

Converting Multi-connected trees into Singly connected trees

nobchap16 - GSview

File Edit Options View Orientation Media Help

File: nobchap16 486, 509pt Page: "367" 19 of 44

16.5. Approximate Inference in Bayesian Networks 367

(a)

C	P(S)
T	.10
F	.90

C	P(R)
T	.80
F	.20

S	R	P(W)
T	T	.99
T	F	.90
F	T	.90
F	F	.00

(b)

C	P(S+R=x)			
C	TT	TF	FT	FF
T	.08	.02	.72	.18
F	.10	.40	.10	.40

S+R	P(W)
T T	.99
T F	.90
F T	.90
F F	.00

Figure 16.11 (a) A multiply connected network with conditional probability tables. (b) A clustered equivalent of the multiply connected network.

Conversion will take exponential time
-Still worth doing if conversion is done off-line
and the cost is amortized over many potential queries

Summary of BN Inference Algorithms

TONS OF APPROACHES

Exact Inference

- Complexity
 - NP-hard (actually #P-Complete; since we “count” models)
 - Polynomial for “Singly connected” networks (one path between each pair of nodes)
- Algorithms
 - Enumeration
 - Variable elimination
 - Avoids the redundant computations of Enumeration
 - [Many others such as “message passing” algorithms, Constraint-propagation based algorithms etc.]

Approximate Inference

- Complexity
 - NP-Hard for both absolute and relative approximation
- Algorithms
 - Based on Stochastic Simulation
 - Sampling from empty networks
 - Rejection sampling
 - Likelihood weighting
 - [And many more]

Performance of approximation algorithms

Absolute approximation: $|P(X|e) - \hat{P}(X|e)| \leq \epsilon$

Relative approximation: $\frac{|P(X|e) - \hat{P}(X|e)|}{P(X|e)} \leq \epsilon$

Relative \Rightarrow absolute since $0 \leq P \leq 1$ (may be $O(2^{-n})$)

Randomized algorithms may fail with probability at most δ

Polytime approximation: $\text{poly}(n, \epsilon^{-1}, \log \delta^{-1})$

Theorem (Dagum and Luby, 1993): both absolute and relative approximation for either deterministic or randomized algorithms are NP-hard for any $\epsilon, \delta < 0.5$

(Absolute approximation polytime with no evidence—Chernoff bounds)

Inference by enumeration

Slightly intelligent way to sum out variables from the joint without actually constructing its explicit representation

Simple query on the burglary network:

$$\begin{aligned}
 & P(B|J = \text{true}, M = \text{true}) \\
 &= P(B, J = \text{true}, M = \text{true}) / P(J = \text{true}, M = \text{true}) \\
 &= \alpha P(B, J = \text{true}, M = \text{true}) \\
 &= \alpha \sum_e \sum_a P(B, e, a, J = \text{true}, M = \text{true})
 \end{aligned}$$

Rewrite full joint entries using product of CPT entries:

$$\begin{aligned}
 & P(B = \text{true} | J = \text{true}, M = \text{true}) \\
 &= \alpha \sum_e \sum_a P(B = \text{true}) P(e) P(a | B = \text{true}, e) P(J = \text{true} | a) P(M = \text{true} | a) \\
 &= \alpha P(B = \text{true}) \sum_e P(e) \sum_a P(a | B = \text{true}, e) P(J = \text{true} | a) P(M = \text{true} | a)
 \end{aligned}$$

Enumeration algorithm

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

ENUMERATION-ASK(X, e, bn) returns a distribution over X

inputs: X , the query variable

e , evidence specified as an event

bn , a belief network specifying joint distribution $P(X_1, \dots, X_n)$

$Q(x) \leftarrow$ a distribution over X

for each value x_i of X do

 extend e with value x_i for X

$Q(x_i) \leftarrow$ ENUMERATE-ALL(VARS(bn), e)

return NORMALIZE($Q(X)$)

ENUMERATE-ALL($vars, e$) returns a real number

if EMPTY?($vars$) then return 1.0

else do

$Y \leftarrow$ FIRST($vars$)

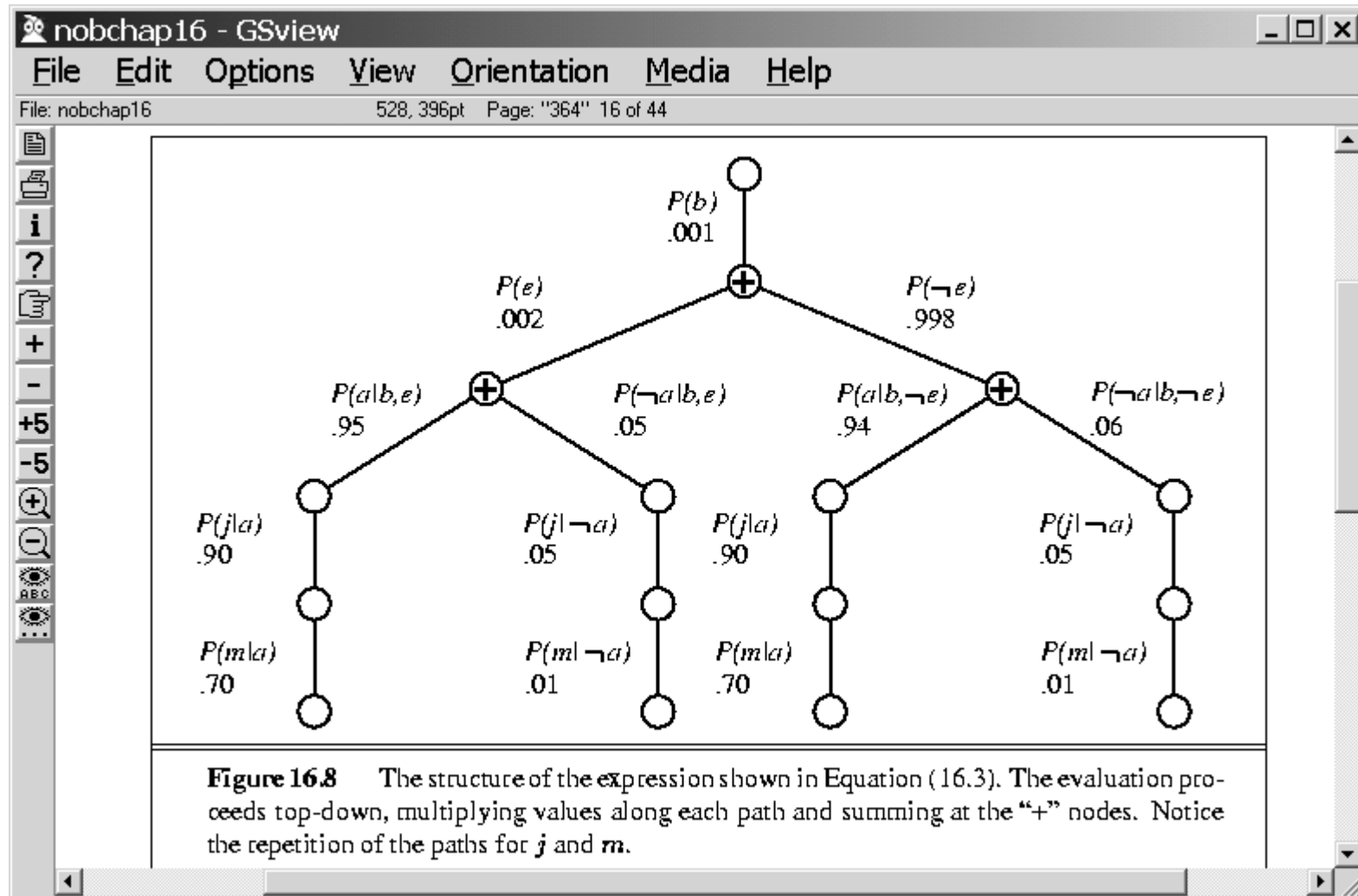
 if Y has value y in e

 then return $P(y \mid Pa(Y)) \times$ ENUMERATE-ALL(REST($vars$), e)

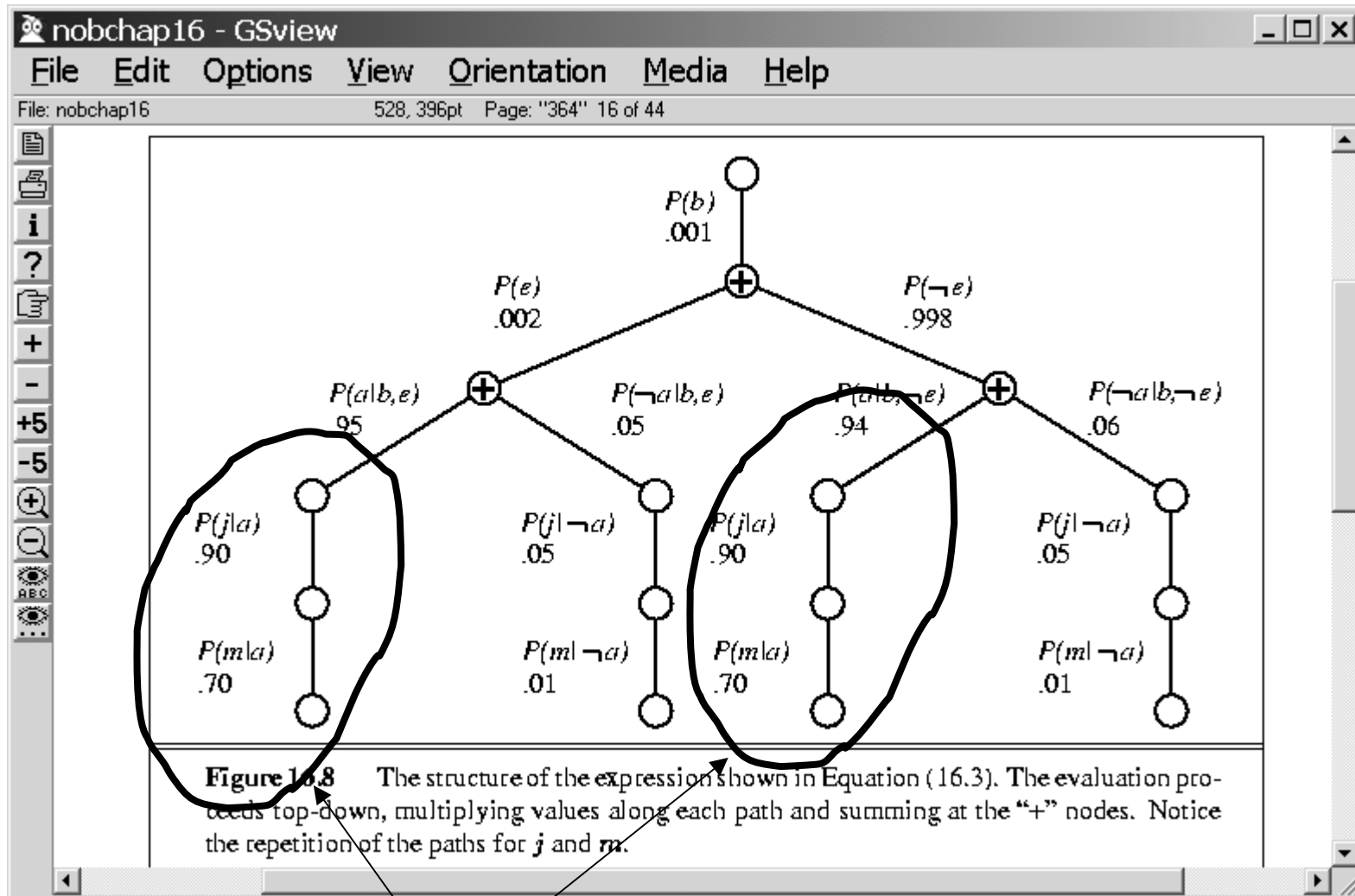
 else return $\sum_y P(y \mid Pa(Y)) \times$ ENUMERATE-ALL(REST($vars$), e_y)

 where e_y is e extended with $Y = y$

Inefficient (redundant) computations in Enumeration



Inefficient (redundant) computations in Enumeration



Repeated multiplications

Inference by variable elimination

Enumeration is inefficient: repeated computation

e.g., computes $P(J = true|a)P(M = true|a)$ for each value of e

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

$$\begin{aligned}
 P(B|J = true, M = true) &= \alpha \underbrace{P(B)}_B \underbrace{\sum_e P(e)}_E \underbrace{\sum_a P(a|B, e)}_A \underbrace{P(J = true|a)}_J \underbrace{P(M = true|a)}_M \\
 &= \alpha P(B) \sum_e P(e) \sum_a P(a|B, e) P(J = true|a) f_M(a) \\
 &= \alpha P(B) \sum_e P(e) \sum_a P(a|B, e) f_I(a) f_M(a) \\
 &= \alpha P(B) \sum_e P(e) \sum_a f_A(a, b, e) f_I(a) f_M(a) \\
 &= \alpha P(B) \sum_e P(e) f_{\bar{A}IM}(b, e) \text{ (sum out } A) \\
 &= \alpha P(B) f_{\bar{E}\bar{A}IM}(b) \text{ (sum out } E) \\
 &= \alpha f_B(b) \times f_{\bar{E}\bar{A}IM}(b)
 \end{aligned}$$

Variable elimination: Basic operations

Pointwise product of factors f_1 and f_2 :

$$f_1(x_1, \dots, x_j, y_1, \dots, y_k) \times f_2(y_1, \dots, y_k, z_1, \dots, z_l) \\ = f(x_1, \dots, x_j, y_1, \dots, y_k, z_1, \dots, z_l)$$

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

Summing out a variable from a product of factors: move any constant factors outside the summation:

$$\sum_x f_1 \times \dots \times f_k = f_1 \times \dots \times f_i \sum_x f_{i+1} \times \dots \times f_k = f_1 \times \dots \times f_i \times f_{\bar{x}}$$

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

Variable elimination algorithm

```
function ELIMINATIONASK( $X, e, \beta_{\pi}$ ) returns a distribution over  $X$ 
inputs:  $X$ , the query variable
        $e$ , evidence specified as an event
        $\beta_{\pi}$ , a belief network specifying joint distribution  $P(X_1, \dots, X_n)$ 
if  $X \in e$  then return observed point distribution for  $X$ 
 $factors \leftarrow []$ ;  $vars \leftarrow REVERSE(VARS[\beta_{\pi}])$ 
for each  $var$  in  $vars$  do
     $factors \leftarrow [MAKEFACTOR(var, e) | factors]$ 
    if  $var$  is a hidden variable then  $factors \leftarrow SUMOUT(var, factors)$ 
return NORMALIZE(POINTWISEPRODUCT( $factors$ ))
```

Inference by stochastic simulation

Basic idea:

- 1) Draw N samples from a sampling distribution S
- 2) Compute an approximate posterior probability \hat{P}
- 3) 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
- MCMC: sample from a stochastic process whose stationary distribution is the true posterior

Sampling from an empty network

```

function PRIORITYSAMPLE(bn) returns an event sampled from  $P(X_1, \dots, X_n)$  specified by bn
  x ← an event with n elements
  for i = 1 to n do
    zi ← a random sample from  $P(X_i \mid \text{Parents}(X_i))$ 
  return x

```

$P(\text{Cloudy}) = \langle 0.5, 0.5 \rangle$

sample → *true*

$P(\text{Sprinkler} \mid \text{Cloudy}) = \langle 0.1, 0.9 \rangle$

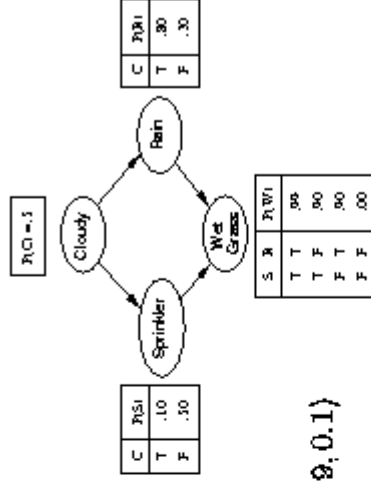
sample → *false*

$P(\text{Rain} \mid \text{Cloudy}) = \langle 0.8, 0.2 \rangle$

sample → *true*

$P(\text{WetGrass} \mid \neg \text{Sprinkler}, \text{Rain}) = \langle 0.9, 0.1 \rangle$

sample → *true*



Sampling from an empty network contd.

Probability that `PRIORSAMPLE` generates a particular event

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

i.e., the true prior probability

Let $N_{PS}(\mathbf{Y} = \mathbf{y})$ be the number of samples generated for which $\mathbf{Y} = \mathbf{y}$, for any set of variables \mathbf{Y} .

Then $\hat{P}(\mathbf{Y} = \mathbf{y}) = N_{PS}(\mathbf{Y} = \mathbf{y})/N$ and

$$\begin{aligned} \lim_{N \rightarrow \infty} \hat{P}(\mathbf{Y} = \mathbf{y}) &= \sum_{\mathbf{h}} S_{PS}(\mathbf{Y} = \mathbf{y}, \mathbf{H} = \mathbf{h}) \\ &= \sum_{\mathbf{h}} P(\mathbf{Y} = \mathbf{y}, \mathbf{H} = \mathbf{h}) \\ &= P(\mathbf{Y} = \mathbf{y}) \end{aligned}$$

That is, estimates derived from `PRIORSAMPLE` are consistent

Rejection sampling

$\hat{P}(X|e)$ estimated from samples agreeing with e

```

function REJECTION SAMPLING( $X, e, b_N, N$ ) returns an approximation to  $P(X|e)$ 
 $N[X] \leftarrow$  a vector of counts over  $X$ , initially zero
for  $j = 1$  to  $N$  do
   $x \leftarrow$  PRIORSAMPLING( $b_N$ )
  if  $x$  is consistent with  $e$  then
     $N[j] \leftarrow N[j] + 1$  where  $z$  is the value of  $X$  in  $x$ 
return NORMALIZE( $N[X]$ )

```

E.g., estimate $P(\text{Rain}|\text{Sprinkler} = \text{true})$ using 100 samples
27 samples have $\text{Sprinkler} = \text{true}$

Of these, 8 have $\text{Rain} = \text{true}$ and 19 have $\text{Rain} = \text{false}$.

$\hat{P}(\text{Rain}|\text{Sprinkler} = \text{true}) = \text{NORMALIZE}(\{8, 19\}) = (0.296, 0.704)$

Similar to a basic real-world empirical estimation procedure

Analysis of rejection sampling

$$\begin{aligned}\hat{P}(X|e) &= \alpha N_{PS}(X, e) && \text{(algorithm defn.)} \\ &= N_{PS}(X, e) / N_{PS}(e) && \text{(normalized by } N_{PS}(e)) \\ &\approx P(X, e) / P(e) && \text{(property of } P_{\text{PRIOR SAMPLE}}) \\ &= P(X|e) && \text{(defn. of conditional probability)}\end{aligned}$$

Hence rejection sampling returns consistent posterior estimates

Problem: hopelessly expensive if $P(e)$ is small

Likelihood weighting

Idea: fix evidence variables, sample only nonevidence variables, and weight each sample by the likelihood it accords the evidence

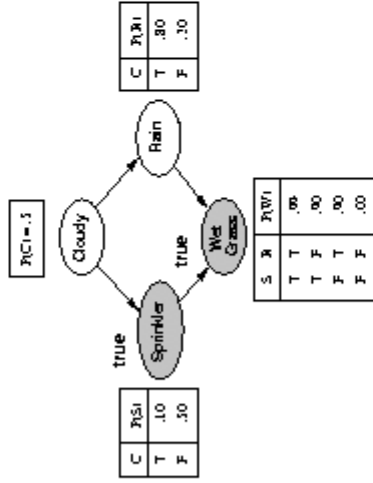
```

function WEIGHTEDSAMPLE( $\beta, e$ ) returns an event and a weight
 $x \leftarrow$  an event with  $n$  elements;  $w \leftarrow 1$ 
for  $i = 1$  to  $n$  do
  if  $X_i$  has a value  $z_i$  in  $e$ 
    then  $w \leftarrow w \times P(X_i = z_i \mid \text{Parents}(X_i))$ 
    else  $z_i \leftarrow$  a random sample from  $P(X_i \mid \text{Parents}(X_i))$ 
  return  $x, w$ 
function LIKELIHOODWEIGHTING( $X, e, \beta, N$ ) returns an approximation to  $P(X \mid e)$ 
 $W[X] \leftarrow$  a vector of weighted counts over  $X$ , initially zero
for  $j = 1$  to  $N$  do
   $x, w \leftarrow$  WEIGHTEDSAMPLE( $\beta, e$ )
   $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
return NORMALIZE( $W[X]$ )

```

Likelihood weighting example

Estimate $P(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$



LW example contd.

Sample generation process:

1. $w \leftarrow 1.0$
2. Sample $P(Cloudy) = (0.5, 0.5)$; say *true*
3. *Sprinkler* has value *true*, so
 $w \leftarrow w \times P(Sprinkler = true | Cloudy = true) = 0.1$
4. Sample $P(Rain | Cloudy = true) = (0.8, 0.2)$; say *true*
5. *WetGrass* has value *true*, so
 $w \leftarrow w \times P(WetGrass = true | Sprinkler = true, Rain = true) = 0.099$

Performance of approximation algorithms

Absolute approximation: $|P(X|e) - \hat{P}(X|e)| \leq \epsilon$

Relative approximation: $\frac{|P(X|e) - \hat{P}(X|e)|}{P(X|e)} \leq \epsilon$

Relative \Rightarrow absolute since $0 \leq P \leq 1$ (may be $O(2^{-n})$)

Randomized algorithms may fail with probability at most δ

Polytime approximation: $\text{poly}(n, \epsilon^{-1}, \log \delta^{-1})$

Theorem (Dagum and Luby, 1993): both absolute and relative approximation for either deterministic or randomized algorithms are NP-hard for any $\epsilon, \delta < 0.5$

(Absolute approximation polytime with no evidence—Chernoff bounds)

Summary of BN Inference Algorithms

TONS OF APPROACHES

Exact Inference

- Complexity
 - NP-hard (actually #P-Complete; since we “count” models)
 - Polynomial for “Singly connected” networks (one path between each pair of nodes)
- Algorithms
 - Enumeration
 - Variable elimination
 - Avoids the redundant computations of Enumeration
 - [Many others such as “message passing” algorithms, Constraint-propagation based algorithms etc.]

Approximate Inference

- Complexity
 - NP-Hard for both absolute and relative approximation
- Algorithms
 - Based on Stochastic Simulation
 - Sampling from empty networks
 - Rejection sampling
 - Likelihood weighting
 - [And many more]