BAYESIAN NETWORKS

AIMA2E CHAPTER 14.1-5 (SOME TOPICS EXCLUDED)

Outline

- ♦ Syntax
- ♦ Semantics
- Parameterized distributions
- ♦ Inference
- Exact inference (enumeration, variable elimination)
- Approximate inference (stochastic simulation)

Bayesian networks

tions A simple, graphical notation for conditional independence asser-

and hence for compact specification of full joint distributions

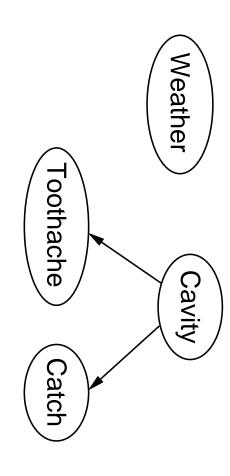
Syntax:

parents Probabilities: a conditional distribution for each node given its **Topology:** a directed, acyclic graph (link \approx "directly influences") Random Variables: a set of nodes, one per variable

$$\mathbf{P}(X_i|Parents(X_i))$$

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

Topology of network encodes conditional independence assertions:

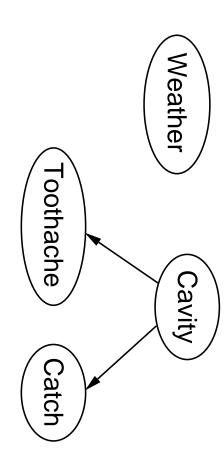


Weather is independent of the other variables

Toothache and Catch are conditionally independent given Cavity

P(Toothache, catch, Cavity, Weather) = ?

Topology of network encodes conditional independence assertions:



Weather is independent of the other variables

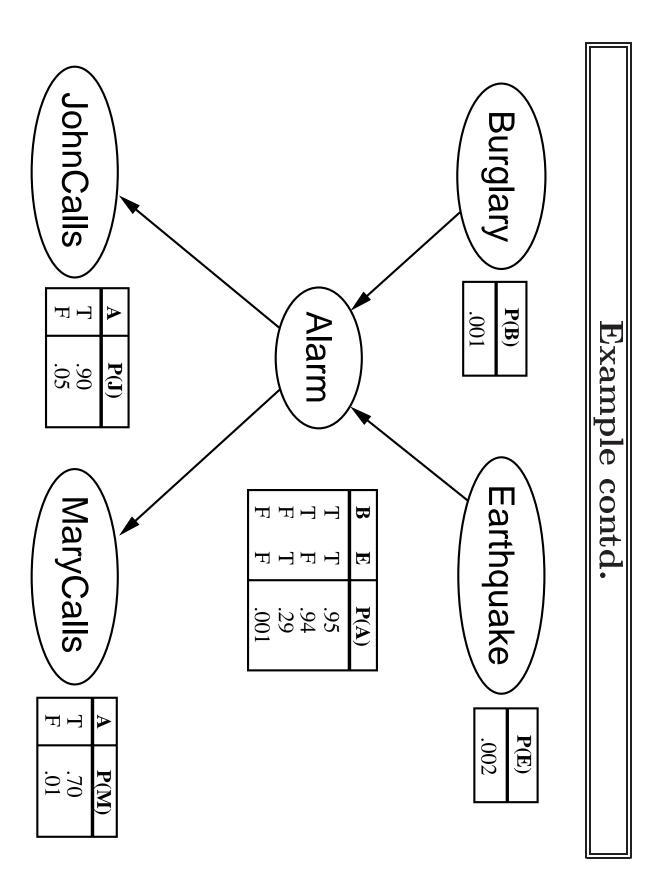
Toothache and Catch are conditionally independent given Cavity

P(Toothache, Catch, Cavity, Weather) = P(Toothache|Cavity)P(Catch|Cavity)P(Cavity)P(Weather)

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

Variables: Burglar, Earthquake, Alarm, JohnCalls, MaryCalls Network topology reflects "causal" knowledge:

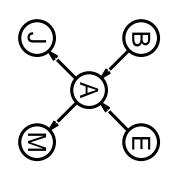
- A burglar can set the alarm off
- An earthquake can set the alarm off
- The alarm can cause Mary to call
- The alarm can cause John to call



${f Compactness}$

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

(the number for $X_i = false$ is just 1-p) Each row requires one number p for $X_i = true$



the complete network requires $O(n \cdot 2^{\kappa})$ numbers If each variable has no more than k parents,

I.e., grows linearly with n, vs. $O(2^n)$ for the full joint distribution

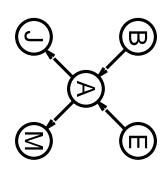
For burglary net, 1+1+4+2+2=10 numbers (vs. $2^5-1=31$)

Global semantics

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

$$\mathbf{P}(X_1,\ldots,X_n) = \prod_{i=1}^n \mathbf{P}(X_i|Parents(X_i))$$

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



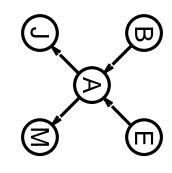
Global semantics

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

$$\mathbf{P}(X_1,\ldots,X_n) = \prod_{i=1}^n \mathbf{P}(X_i|Parents(X_i))$$

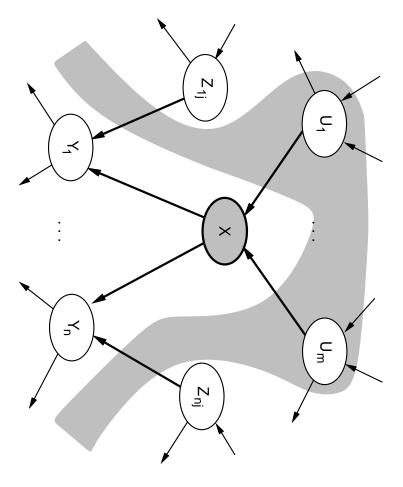
e.g.,
$$P(j \land m \land a \land \neg b \land \neg e)$$

$$= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e)$$



Local semantics

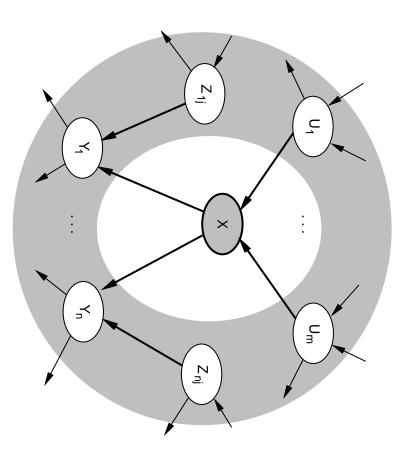
of its nondescendants given its parents Local semantics: each node is conditionally independent



Theorem: Local semantics \Leftrightarrow global semantics

Markov blanket

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



Constructing Bayesian networks

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

- 1. Choose an ordering of variables X_1, \ldots, X_n
- 2. For i = 1 to nadd X_i to the network select parents from X_1, \ldots, X_{i-1} such that $\mathbf{P}(X_i|Parents(X_i)) = \mathbf{P}(X_i|X_1, \dots, X_{i-1})$

This choice of parents guarantees the global semantics:

$$\mathbf{P}(X_1,\ldots,X_n) = \prod_{i=1}^n \mathbf{P}(X_i|X_1,\ldots,X_{i-1}) \quad \text{(chain rule)}$$
$$= \prod_{i=1}^n \mathbf{P}(X_i|Parents(X_i)) \quad \text{(by construction)}$$

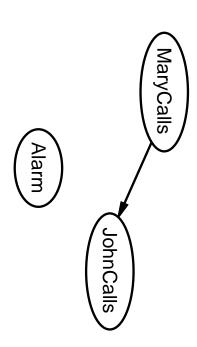
Suppose we choose the ordering M, J, A, B, E





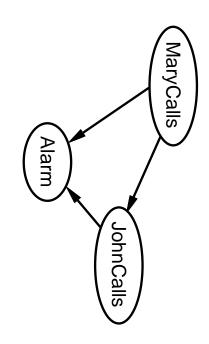
$$P(J|M) = P(J)$$
?

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



$$P(J|M) = P(J)$$
? No $P(A|J,M) = P(A|J)$?

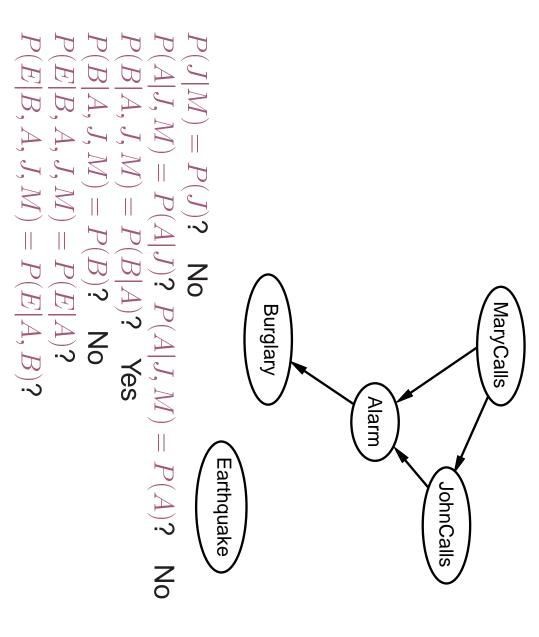
Suppose we choose the ordering M, J, A, B, E



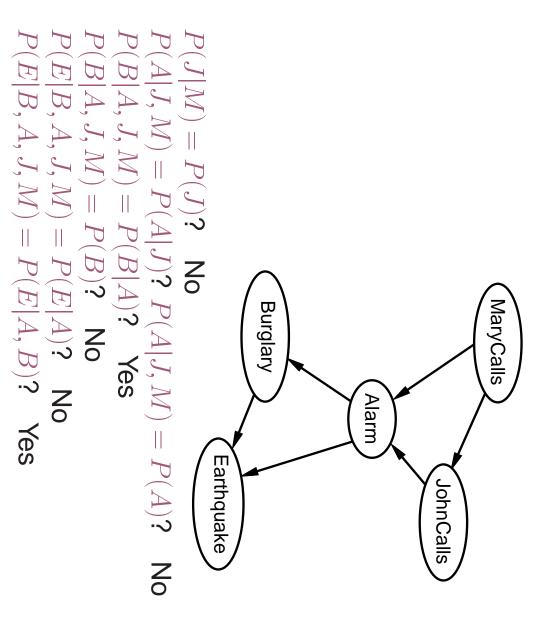
Burglary

$$P(J|M) = P(J)$$
? No $P(A|J,M) = P(A|J)$? P($A|J,M$) = $P(A|J)$? P($A|J,M$) = $P(B|A,J,M) = P(B|A)$? No $P(B|A,J,M) = P(B)$?

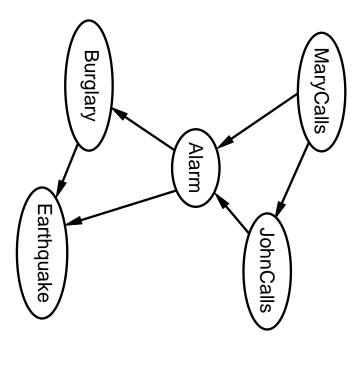
Suppose we choose the ordering M, J, A, B, E



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



Example contd.



Deciding conditional independence is hard in noncausal directions

humans!) (Causal models and conditional independence seem hardwired for

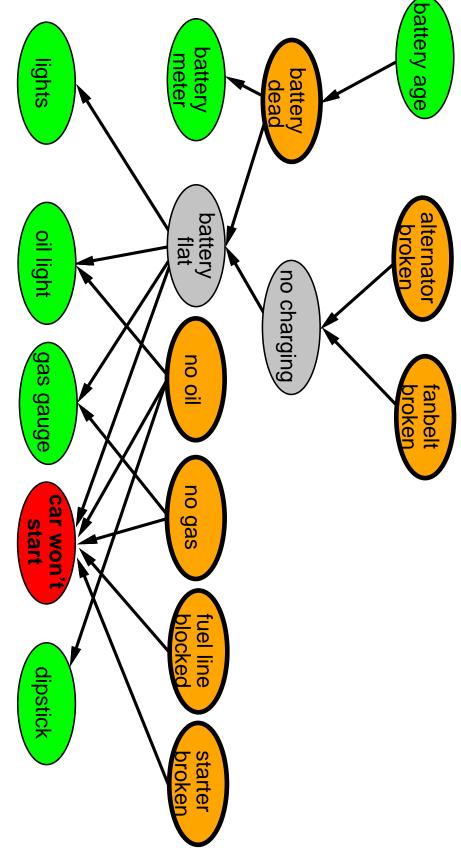
Assessing conditional probabilities is hard in noncausal directions

Network is less compact: 1+2+4+2+4=13 numbers needed

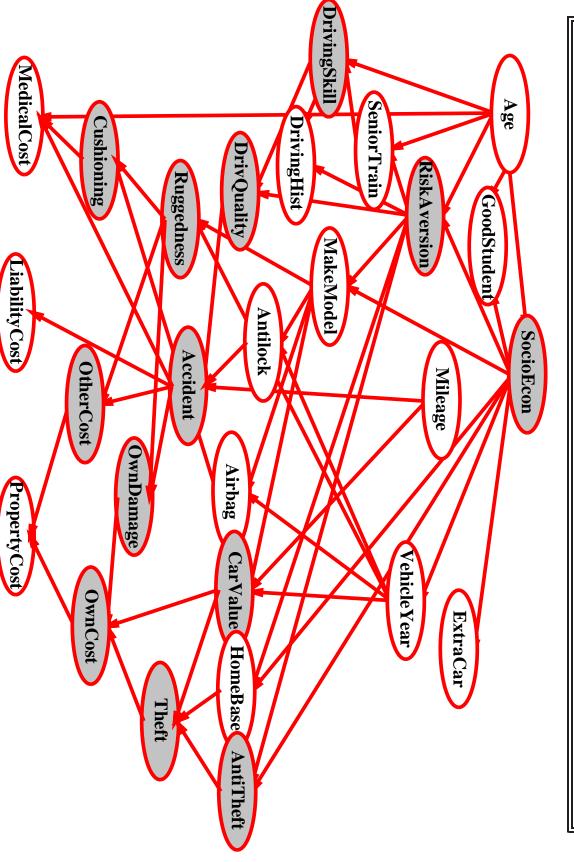
Example: Car diagnosis

Initial evidence: car won't start

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



Example: Car insurance



Compact conditional distributions

CPT becomes infinite with continuous-valued parent or child CPT grows exponentially with no. of parents

Solution: canonical distributions that are defined compactly

Deterministic nodes are the simplest case:

$$X = f(Parents(X))$$
 for some function f

E.g., Boolean functions

$$NorthAmerican \Leftrightarrow Canadian \lor US \lor Mexican$$

E.g., numerical relationships among continuous variables

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

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

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

0.02	T F F 0.4 0.6	0 06	0.12
	$0.02 = 0.2 \times 0.1$		$0.02 = 0.2 \times 0.1$ 0.6 $0.06 = 0.6 \times 0.1$ $0.12 = 0.6 \times 0.2$

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

Continuous nodes

Networks may have discrete RVs, continuous RVs, or a mix of the

All continuous (e.g., conditional Gaussian). Linear dynamic sysyems (Kalman filter).

Gaussian mixture models Discrete parents, continuous children (e.g., conditional Gaussian).

tions). Difficult to deal with. Continuous parents, discrete children (e.g., probit and logit func-

Inference tasks

Simple queries: compute posterior marginal $P(X_i|\mathbf{E}=\mathbf{e})$ $\textbf{e.g.,}\ P(NoGas|Gauge = empty, Lights = on, Starts = false)$

Conjunctive queries: $P(X_i, X_j | \mathbf{E} = \mathbf{e}) = P(X_i | \mathbf{E} = \mathbf{e})P(X_j | X_i, \mathbf{E} = \mathbf{e})$

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?

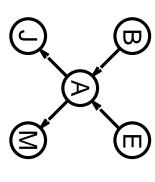
Explanation: why do I need a new starter motor?

Inference by enumeration

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

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:

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

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

Enumeration algorithm

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

inputs: X, the query variable

e, observed values for variables E

bn, a Bayesian network with variables $\{X\} \cup \mathbf{E} \cup \mathbf{Y}$

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

for each value x_i of X do

extend e with value x_i for X

 $\mathbf{Q}(x_i) \leftarrow \text{Enumerate-All}(\text{Vars}[bn], \mathbf{e})$

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

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

if EMPTY?(vars) then return 1.0

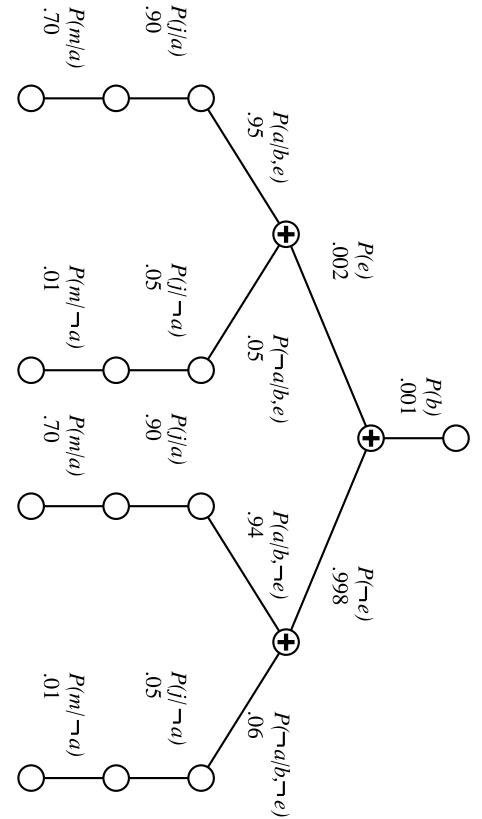
 $Y \leftarrow \text{FIRST}(vars)$

if Y has value y in e

else return Σ_y $P(y \mid Pa(Y)) \times \text{ENUMERATE-ALL(REST(<math>vars$), \mathbf{e}_y) then return $P(y \mid Pa(Y)) \times \text{Enumerate-All(Rest(vars), e)}$ where \mathbf{e}_y is \mathbf{e} extended with Y = y

Evaluation tree

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



Inference by variable elimination

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

$$\begin{split} \mathbf{P}(B|j,m) \\ &= \alpha \underbrace{\mathbf{P}(B)}_{B} \sum_{e} \underbrace{P(e)}_{E} \sum_{a} \underbrace{\mathbf{P}(a|B,e)}_{A} \underbrace{P(j|a)}_{A} \underbrace{P(m|a)}_{M} \\ &= \alpha \mathbf{P}(B) \sum_{e} P(e) \sum_{a} \mathbf{P}(a|B,e) P(j|a) f_{M}(a) \\ &= \alpha \mathbf{P}(B) \sum_{e} P(e) \sum_{a} \mathbf{P}(a|B,e) f_{J}(a) f_{M}(a) \\ &= \alpha \mathbf{P}(B) \sum_{e} P(e) \sum_{a} f_{A}(a,b,e) f_{J}(a) f_{M}(a) \\ &= \alpha \mathbf{P}(B) \sum_{e} P(e) f_{\bar{A}JM}(b,e) \text{ (sum out } A) \\ &= \alpha f_{B}(b) \times f_{\bar{E}\bar{A}JM}(b) \text{ (sum out } E) \\ &= \alpha f_{B}(b) \times f_{\bar{E}\bar{A}JM}(b) \text{ (sum out } E) \end{split}$$

Variable elimination: Basic operations

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

$$\sum_x f_1 imes \cdots imes f_i imes f_i imes f_{i+1} imes \cdots imes f_k = f_1 imes \cdots imes f_i imes f_{ar{X}}$$

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

Pointwise product of factors f_1 and f_2 :

$$f_1(x_1, ..., x_j, y_1, ..., y_k) \times f_2(y_1, ..., y_k, z_1, ..., z_l)$$

= $f(x_1, ..., x_j, y_1, ..., y_k, z_1, ..., z_l)$

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

Variable elimination algorithm

```
function ELIMINATION-ASK(X, \mathbf{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 $\mathbf{P}(X_1,\ldots,X_n)$

 $factors \leftarrow []; vars \leftarrow Reverse(Vars[bn])$

for each var in vars do

 $factors \leftarrow [Make-Factor(var, \mathbf{e})|factors]$

if var is a hidden variable then $factors \leftarrow Sum-Out(var, factors)$

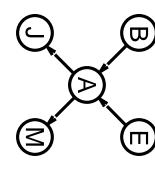
return NORMALIZE(POINTWISE-PRODUCT(factors))

Irrelevant variables

Consider the query P(JohnCalls|Burglary = true)

$$P(J|b) = \alpha P(b) \sum_{e} P(e) \sum_{a} P(a|b,e) P(J|a) \sum_{m} P(m|a)$$

Sum over m is identically 1; M is **irrelevant** to the query



Thm 1: Y is irrelevant unless $Y \in Ancestors(\{X\} \cup \mathbb{E})$

so M is irrelevant $Ancestors(\{X\} \cup \mathbf{E}) = \{Alarm, Earthquake\}$ Here, X = JohnCalls, $\mathbb{E} = \{Burglary\}$, and

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 $O(d^{\kappa}n)$

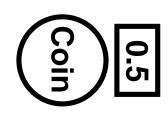
Multiply connected networks:

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

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

Sampling from an empty network

function PRIOR-SAMPLE(bn) **returns** an event sampled from bn**inputs**: bn, a belief network specifying joint distribution $\mathbf{P}(X_1,\ldots,X_n)$

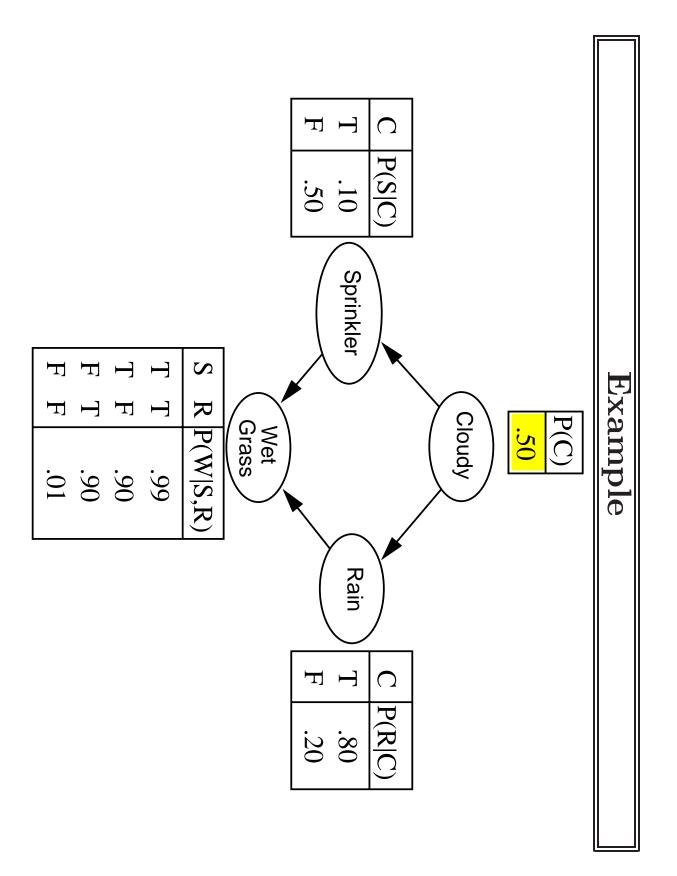
 $\mathbf{x} \leftarrow$ an event with n elements

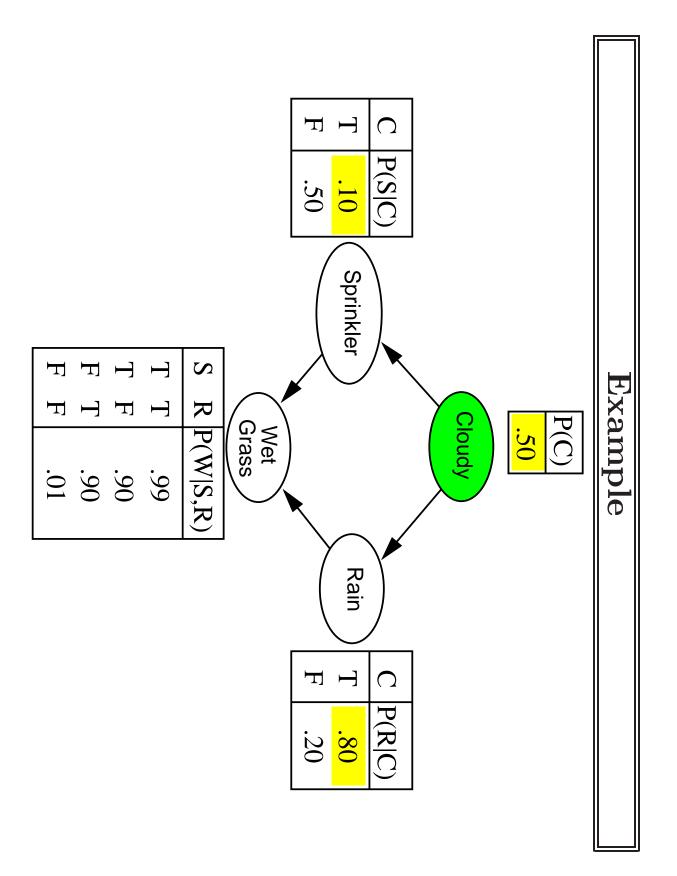
for i = 1 to n do

return x

 $x_i \leftarrow a \text{ random sample from } \mathbf{P}(X_i \mid Parents(X_i))$

AIMA2e Chapter 14.1–5 (some topics excluded)





Sampling from an empty network contd.

Probability that PRIORSAMPLE generates a particular event $S_{PS}(x_1 ... x_n) = \prod_{i=1}^n P(x_i | Parents(X_i)) = P(x_1 ... x_n)$

i.e., the true prior probability

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

 x_1,\ldots,x_n Let $N_{PS}(x_1 \ldots x_n)$ be the number of samples generated for event

Then we have

$$\lim_{N \to \infty} \hat{P}(x_1, \dots, x_n) = \lim_{N \to \infty} N_{PS}(x_1, \dots, x_n)/N$$
$$= S_{PS}(x_1, \dots, x_n)$$
$$= P(x_1 \dots x_n)$$

That is, estimates derived from PRIORSAMPLE are consistent

Shorthand: $P(x_1, ..., x_n) \approx P(x_1 ... x_n)$

Rejection sampling

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

function Rejection-Sampling(X, e, bn, N) returns an estimate of P(X|e)**local variables:** N, a vector of counts over X, initially zero

for j = 1 to N do

 $\mathbf{x} \leftarrow \text{Prior-Sample}(bn)$

if x is consistent with e then

 $N[x] \leftarrow N[x] + 1$ where x is the value of X in x

return Normalize(N[X])

E.g., estimate $\mathbf{P}(Rain|Sprinkler = true)$ using 100 samples 27 samples have Sprinkler = true Of these, 8 have Rain = true and 19 have Rain = false.

 $\hat{\mathbf{P}}(Rain|Sprinkler = true) = \text{Normalize}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle$

Similar to a basic real-world empirical estimation procedure

Analysis of rejection sampling

$$\hat{\mathbf{P}}(X|\mathbf{e}) = \alpha \mathbf{N}_{PS}(X,\mathbf{e})$$
 (algorithm defn.)
 $= \mathbf{N}_{PS}(X,\mathbf{e})/N_{PS}(\mathbf{e})$ (normalized by $N_{PS}(\mathbf{e})$)
 $\approx \mathbf{P}(X,\mathbf{e})/P(\mathbf{e})$ (property of PRIORSAMPLE)
 $= \mathbf{P}(X|\mathbf{e})$ (defn. of conditional probability)

Hence rejection sampling returns consistent posterior estimates

Problem: hopelessly expensive if $P(\mathbf{e})$ is small

 $P(\mathbf{e})$ drops off exponentially with number of evidence variables!

Likelihood weighting

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

function LIKELIHOOD-WEIGHTING(X, e, bn, N) returns an estimate of P(X|e)**local variables: W**, a vector of weighted counts over X, initially zero

for j = 1 to N do

 $\mathbf{x}, w \leftarrow \text{Weighted-Sample}(bn)$

 $\mathbf{W}[x] \leftarrow \mathbf{W}[x] + w$ where x is the value of X in \mathbf{x}

return Normalize(W[X])

function Weighted-Sample(bn, e) returns an event and a weight

 $\mathbf{x} \leftarrow$ an event with n elements; $w \leftarrow$ 1

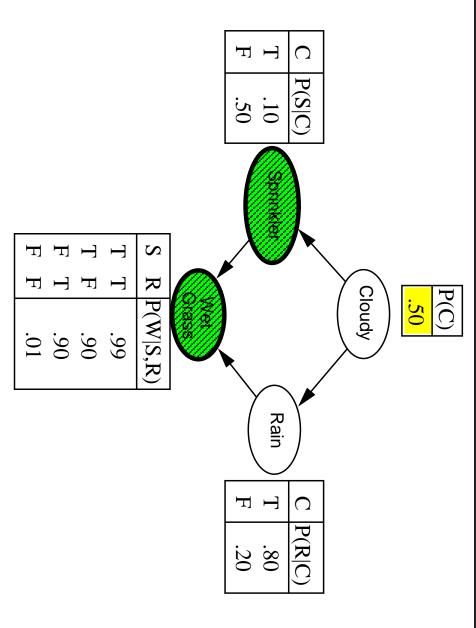
for i = 1 to n do

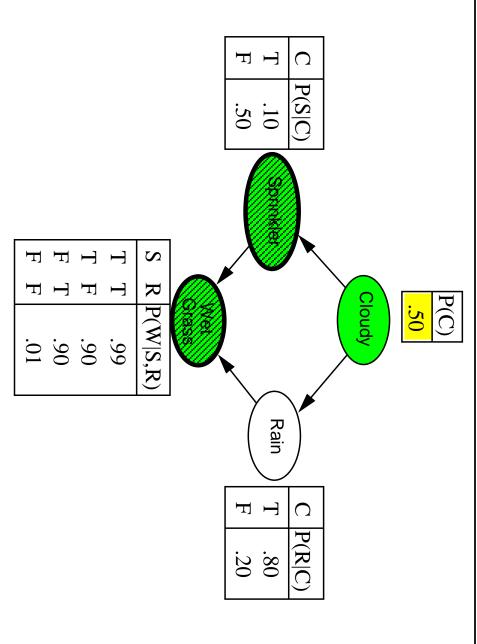
if X_i has a value x_i in **e**

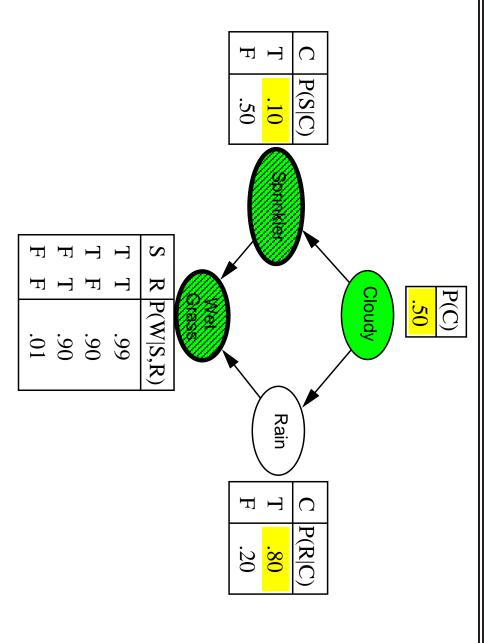
then $w \leftarrow w \times P(X_i = x_i \mid Parents(X_i))$

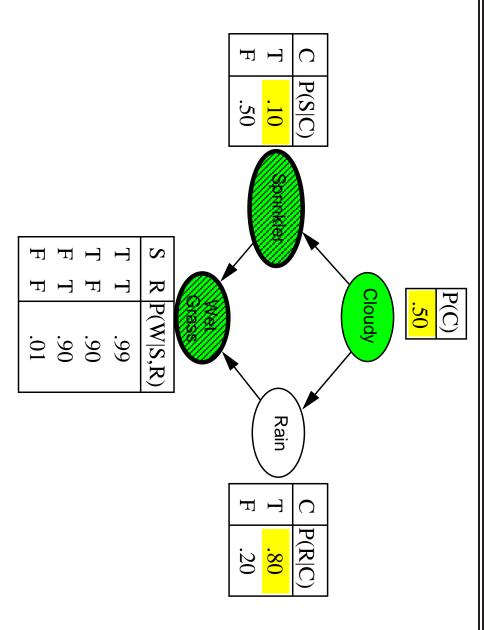
else $x_i \leftarrow$ a random sample from $P(X_i \mid Parents(X_i))$

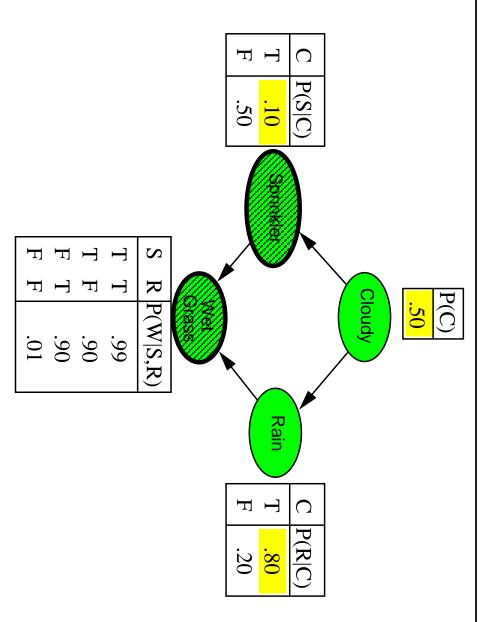
eturn x, w

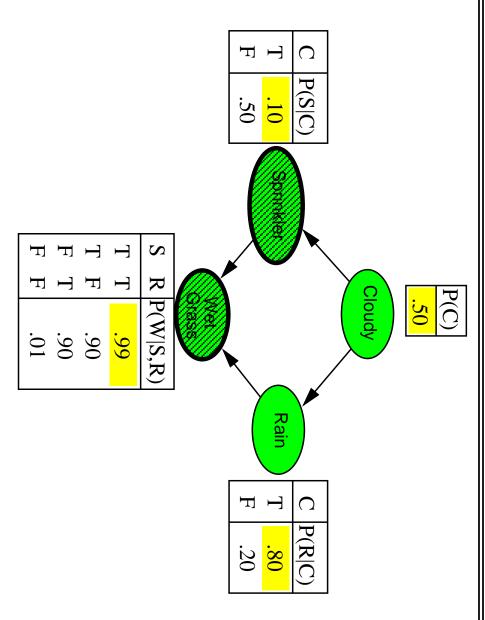


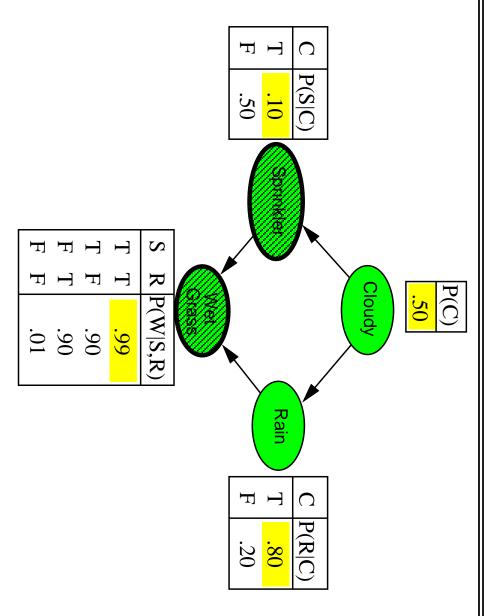












 $w = 1.0 \times 0.1 \times 0.99 = 0.099$

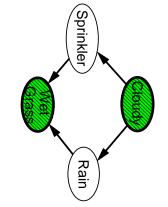
Likelihood weighting analysis

Sampling probability for WEIGHTEDSAMPLE is

$$S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^{t} P(z_i | Parents(Z_i))$$

Note: pays attention to evidence in ancestors only

somewhere "in between" prior and posterior distribution



Weight for a given sample z, e is

$$w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^{m} P(e_i | Parents(E_i))$$

Weighted sampling probability is

$$S_{WS}(\mathbf{z}, \mathbf{e})w(\mathbf{z}, \mathbf{e})$$

$$= \prod_{i=1}^{l} P(z_i|Parents(Z_i)) \quad \prod_{i=1}^{m} P(e_i|Parents(E_i))$$

$$= P(\mathbf{z}, \mathbf{e}) \text{ (by standard global semantics of network)}$$

because a few samples have nearly all the total weight but performance still degrades with many evidence variables Hence likelihood weighting returns consistent estimates