

CSE 573: Artificial Intelligence Autumn 2012

Bayesian Networks

Dan Weld

Many slides adapted from Dan Klein, Stuart Russell, Andrew Moore & Luke Zettlemoyer

Outline

- Probabilistic models (and inference)
 - Bayesian Networks (BNs)
 - Independence in BNs
 - Efficient Inference in BNs
 - Learning
- Whirlwind, so...
 - Take CSE 515 (Statistical Methods)
 - Ben Taskar, Spring 2013

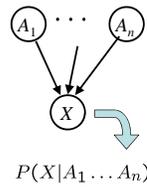
Bayes' Nets: Big Picture

- Two problems with using full joint distribution tables as our probabilistic models:
 - Unless there are only a few variables, the joint is WAY too big to represent explicitly
 - Hard to learn (estimate) anything empirically about more than a few variables at a time
- Bayes' nets: a technique for describing complex joint distributions (models) using simple, local distributions (conditional probabilities)
 - Aka **graphical model**
 - We describe how variables locally interact
 - Local interactions chain together to give global, indirect interactions

Bayes' Net Semantics

Formally:

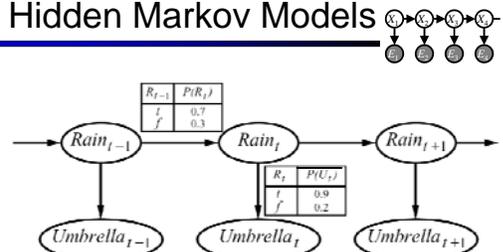
- A set of **nodes**, one per random variable
- **Directed edges**, forming an **acyclic** graph
- A **CPT for each node**
 - CPT = "Conditional Probability Table"
 - Collection of distributions over X, one for each combination of parents' values



$P(X|A_1 \dots A_n)$

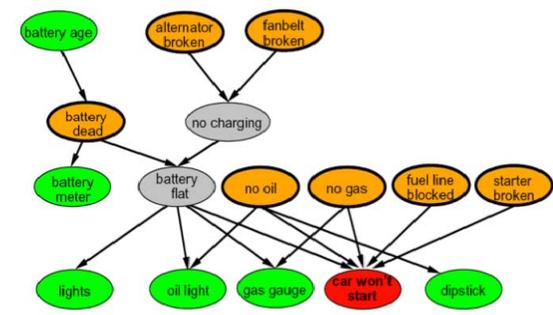
A Bayes Net = Topology (graph) + Local Conditional Probabilities

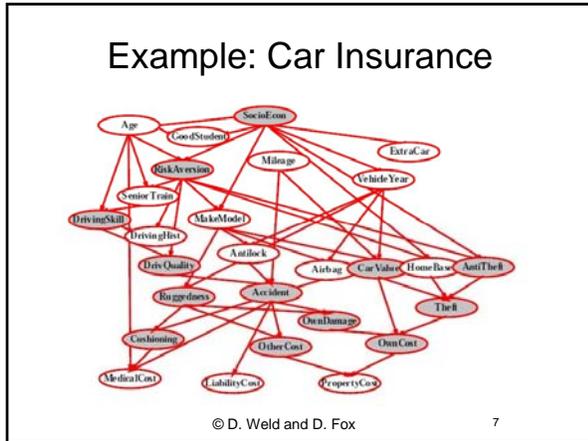
Hidden Markov Models



- An HMM is defined by:
 - Initial distribution: $P(X_1)$
 - Transitions: $P(X_i|X_{i-1})$
 - Emissions: $P(E_i|X_i)$

Example Bayes' Net: Car





Probabilities in BNs

- Bayes' nets **implicitly** encode joint distributions
 - As a product of local conditional distributions
 - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

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

- Does this always work? Why?
- Not every BN can represent every joint distribution
 - The topology enforces certain *independence* assumptions
 - Compare to the exact decomposition according to the chain rule!

Example: Independent Coin Flips

- N independent coin flips

$2^n - 1$

- No interactions between variables $X \perp\!\!\!\perp Y$

Conditional Independence

- Unconditional (absolute) independence very rare (why?)
- Conditional independence* is our most basic and robust form of knowledge about uncertain environments:

$$\forall x, y, z : P(x, y|z) = P(x|z)P(y|z) \quad X \perp\!\!\!\perp Y | Z$$

$$\forall x, y, z : P(x|z, y) = P(x|z)$$
- What about fire, smoke, alarm?

Example: Alarm Network

- Variables**
 - B: Burglary
 - A: Alarm goes off
 - M: Mary calls
 - J: John calls
 - E: Earthquake!
- How big is joint distribution?
 - $2^n - 1 = 31$ parameters

Example: Alarm Network

Only 10 params

B	P(B)
+b	0.001
-b	0.999

E	P(E)
+e	0.002
-e	0.998

A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-e	+a	0.94
+b	-e	-a	0.06
-b	+e	+a	0.29
-b	+e	-a	0.71
-b	-e	+a	0.001
-b	-e	-a	0.999

Example: Traffic II

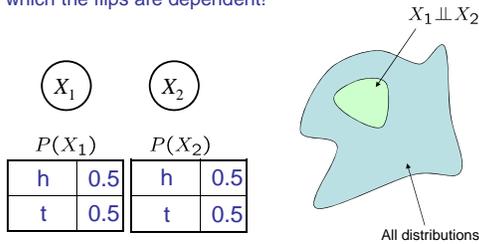
- Let's build a graphical model
- Variables
 - T: Traffic
 - R: It rains
 - L: Low pressure
 - D: Roof drips
 - B: Ballgame
 - C: Cavity

Changing Bayes' Net Structure

- The same joint distribution can be encoded in many different Bayes' nets
- Analysis question: given some edges, what other edges do you need to add?
 - One answer: fully connect the graph
 - Better answer: don't make any false conditional independence assumptions

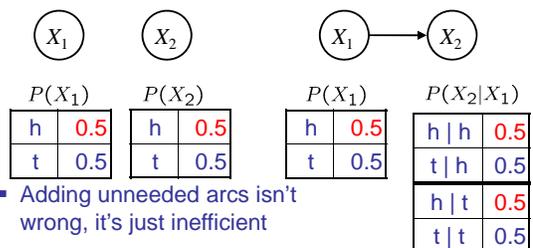
Example: Independence

- For this graph, you can fiddle with (the CPTs) all you want, but you won't be able to represent any distribution in which the flips are dependent!



Example: Coins

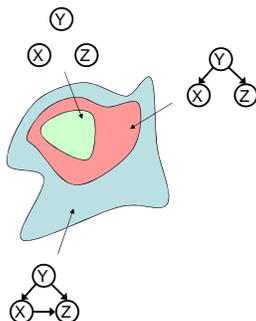
- Extra arcs don't prevent representing independence, just allow non-independence



- Adding unneeded arcs isn't wrong, it's just inefficient

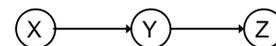
Topology Limits Distributions

- Given some graph topology G , only certain joint distributions can be encoded
- The graph structure guarantees certain (conditional) independences
- (There might be more independence)
- Adding arcs increases the set of distributions, but has several costs
- Full conditioning can encode any distribution



Independence in a BN

- Important question about a BN:
 - Are two nodes independent given certain evidence?
 - If yes, can prove using algebra (tedious in general)
 - If no, can prove with a counter example
 - Example:



- Question: are X and Z independent?
 - Answer: no.
 - Example: low pressure causes rain, which causes traffic.
 - Knowledge about X may change belief in Z,
 - Knowledge about Z may change belief in X (via Y)
 - Addendum: they *could* be independent: how?

Causal Chains

- This configuration is a "causal chain"



X: Low pressure
Y: Rain
Z: Traffic

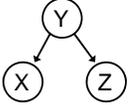
$$P(x, y, z) = P(x)P(y|x)P(z|y)$$
 - Is X independent of Z given Y?

$$P(z|x, y) = \frac{P(x, y, z)}{P(x, y)} = \frac{P(x)P(y|x)P(z|y)}{P(x)P(y|x)} = P(z|y) \quad \text{Yes!}$$

Evidence along the chain "blocks" the influence

Common Parent

- Another basic configuration: two effects of the same parent

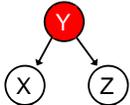


Y: Project due
X: Forum busy
Z: Lab full

 - Are X and Z independent?
 - Are X and Z independent given Y?

Common Parent

- Another basic configuration: two effects of the same parent



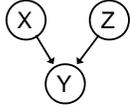
Y: Project due
X: Forum busy
Z: Lab full

 - Are X and Z independent?
 - Are X and Z independent given Y?

$$P(z|x, y) = \frac{P(x, y, z)}{P(x, y)} = \frac{P(y)P(x|y)P(z|y)}{P(y)P(x|y)} = P(z|y) \quad \text{Yes!}$$
 - Observing the cause blocks influence between effects.

Common Effect

- Last configuration: two causes of one effect (v-structures)

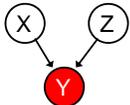


X: Raining
Z: Ballgame
Y: Traffic

 - Are X and Z independent?
 - Yes: the ballgame and the rain cause traffic, but they are not correlated
 - Still need to prove they must be (try it!)

Common Effect

- Last configuration: two causes of one effect (v-structures)



X: Raining
Z: Ballgame
Y: Traffic

 - Are X and Z independent?
 - Yes: the ballgame and the rain cause traffic, but they are not correlated
 - Still need to prove they must be (try it!)
 - Are X and Z independent given Y?
 - No: seeing traffic puts the rain and the ballgame in competition as explanation!
 - This is backwards from the other cases
 - Observing an effect **activates** influence between possible causes.

The General Case

- Any complex example can be analyzed using these three canonical cases
- General question: in a given BN, are two variables independent (given evidence)?
- Solution: analyze the graph

Reachability (D-Separation)

- Question: Are X and Y conditionally independent given evidence vars {Z}?
- Yes, if X and Y "separated" by Z
- Look for active paths from X to Y
- No active paths = independence!
- A path is active if each triple is active:
- Causal chain $A \rightarrow B \rightarrow C$ where B is **unobserved** (either direction)
- Common cause $A \leftarrow B \rightarrow C$ where B is **unobserved**
- Common effect (aka v-structure) $A \rightarrow B \leftarrow C$ where B or one of its descendants is **observed**
- All it takes to block a path is a single inactive segment

Active Triples

Inactive Triples

Example: Independent?

Active Segments

$R \perp\!\!\!\perp B$ Yes

$R \perp\!\!\!\perp B | T$ No

$R \perp\!\!\!\perp B | T'$ No

Example: Independent?

Active Segments

$L \perp\!\!\!\perp T' | T$ Yes

$L \perp\!\!\!\perp B$ Yes

$L \perp\!\!\!\perp B | T$ No

$L \perp\!\!\!\perp B | T'$ No

$L \perp\!\!\!\perp B | T, R$ Yes

Example

Active Segments

- Variables:
 - R: Raining
 - T: Traffic
 - D: Roof drips
 - S: I'm sad
- Questions:

$T \perp\!\!\!\perp D$ No

$T \perp\!\!\!\perp D | R$ Yes

$T \perp\!\!\!\perp D | R, S$ No

Given **Markov Blanket**, X is Independent of All Other Nodes

$MB(X) = \text{Par}(X) \cup \text{Childs}(X) \cup \text{Par}(\text{Childs}(X))$

© D. Weld and D. Fox 43

Given **Markov Blanket**, X is Independent of All Other Nodes

$MB(X) = \text{Par}(X) \cup \text{Childs}(X) \cup \text{Par}(\text{Childs}(X))$

© D. Weld and D. Fox 44

Summary

- Bayes nets compactly encode joint distributions (JDs)
 - Other graphical models too: factor graphs, CRFs, ...
- Guaranteed independencies of distributions can be deduced from BN graph structure
- D-separation gives precise conditional independence guarantees from graph alone
- A Bayes' net's JD may have further (conditional) independence known only from specific CPTs

Outline

- Probabilistic models (and inference)
 - Bayesian Networks (BNs)
 - Independence in BNs
 - Efficient Inference in BNs**
 - Learning

Inference in BNs

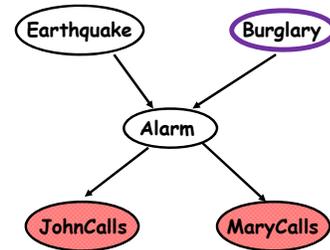
This graphical independence representation yields efficient inference schemes

- We generally want to compute
 - Marginal probability: $Pr(Z)$,
 - $Pr(Z|E)$ where E is (conjunctive) evidence
 - Z : query variable(s),
 - E : evidence variable(s)
 - everything else: hidden variable

Computations organized by network topology

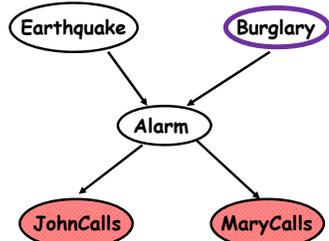
© D. Weld and D. Fox 54

$P(B | J=true, M=true)$



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

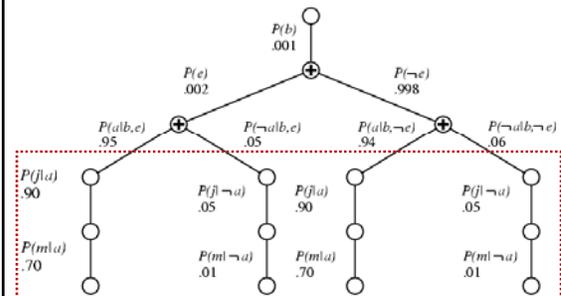
$P(B | J=true, M=true)$



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

Variable Elimination

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

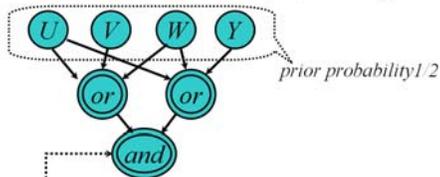


Repeated computations → Dynamic Programming

Reducing 3-SAT to Bayes Nets

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

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



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

© Jack Breese (Microsoft) & Daphne Koller (Stanford)

70

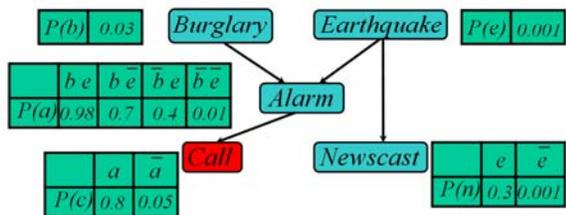
Approximate Inference in Bayes Nets Sampling based methods

(Based on slides by Jack Breese and Daphne Koller)

67

Bayes Net is a generative model

- We can easily generate **samples** from the distribution represented by the Bayes net
- Generate one variable at a time in **topological order**

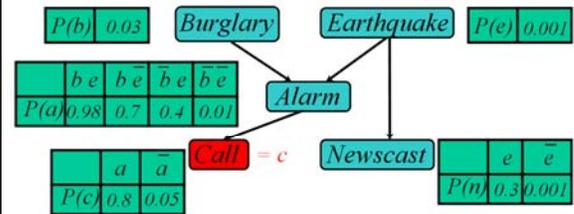


Use the samples to compute probabilities, say $P(c)$ or $P(n|c)$

© Jack Breese (Microsoft) & Daphne Koller (Stanford)

70

Stochastic simulation $P(B|C)$

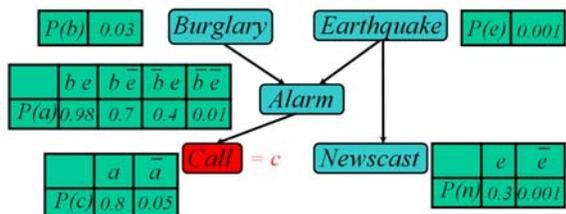


© Jack Breese (Microsoft) & Daphne Koller (Stanford)

69

71

Stochastic simulation $P(B|C)$



Samples:

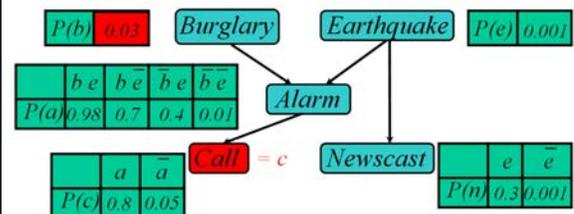
B	E	A	C	N

© Jack Breese (Microsoft) & Daphne Koller (Stanford)

70

71

Stochastic simulation $P(B|C)$



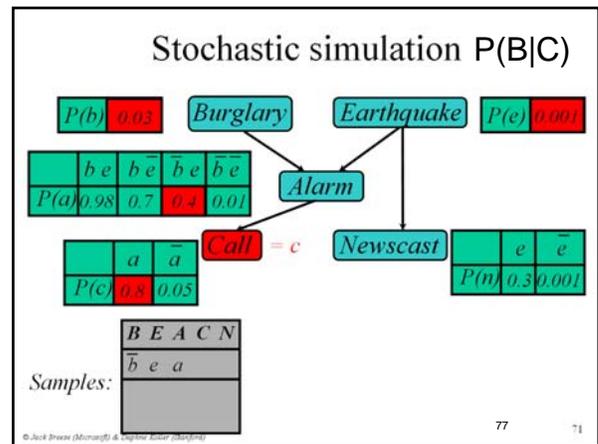
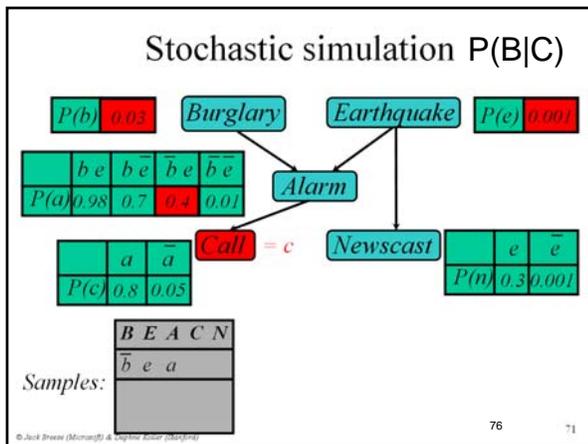
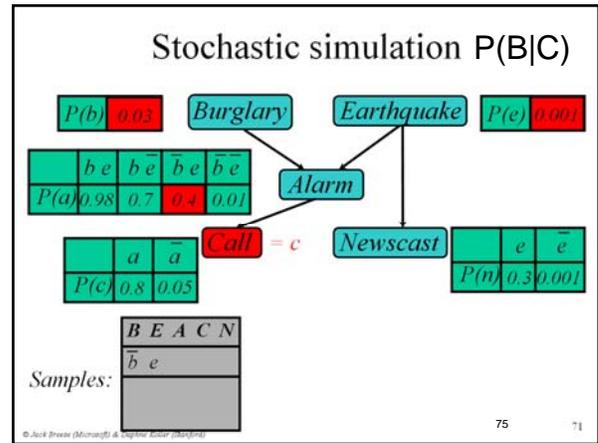
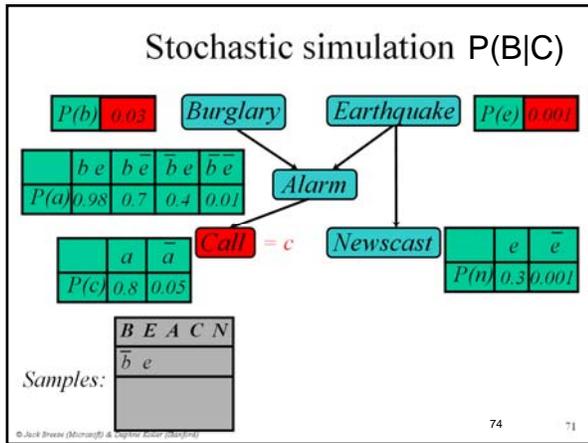
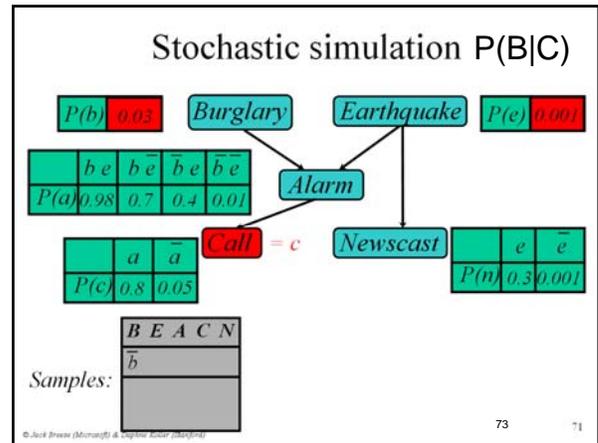
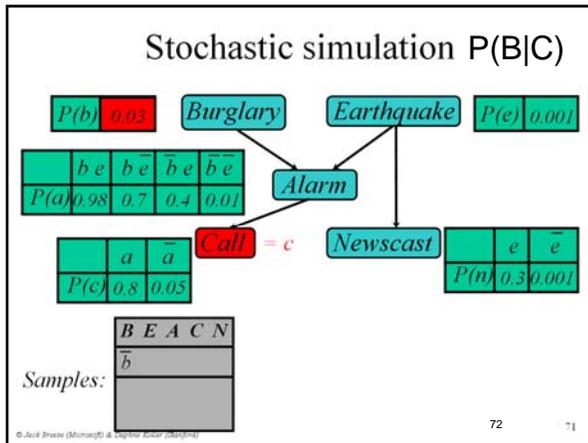
Samples:

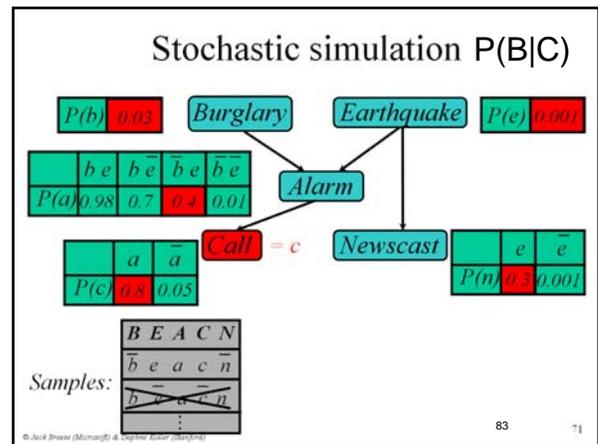
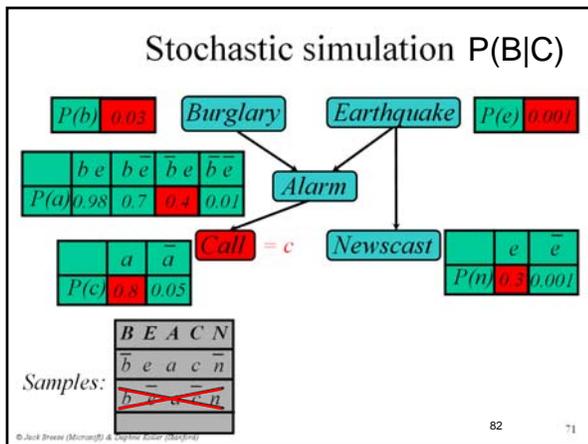
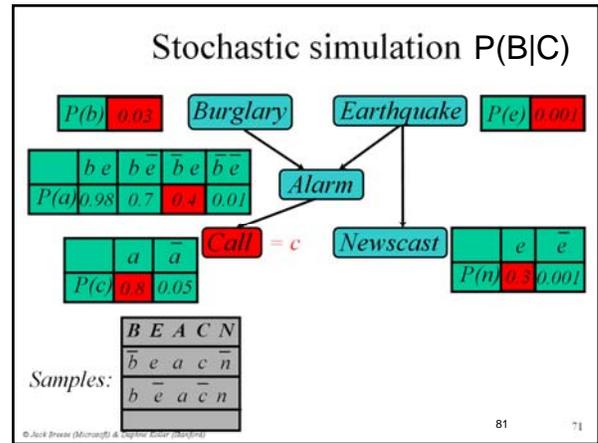
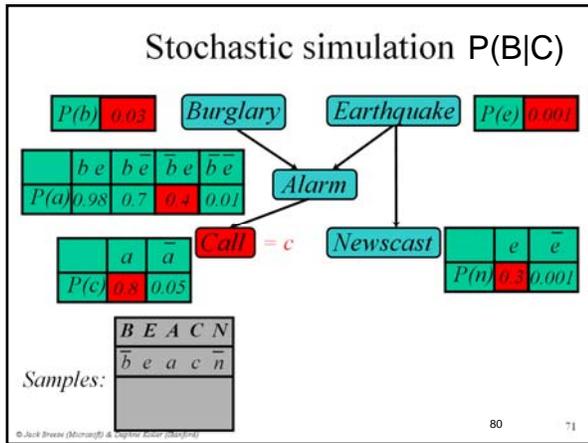
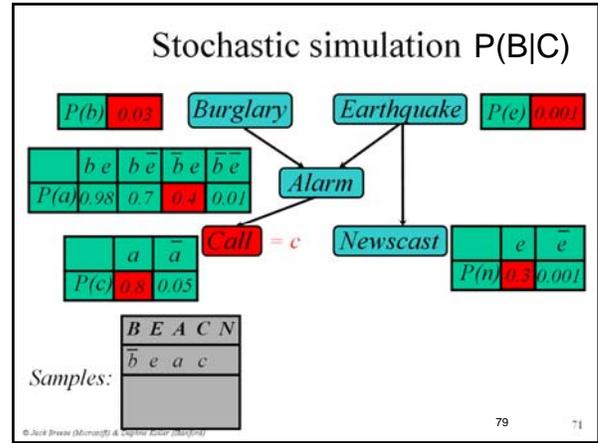
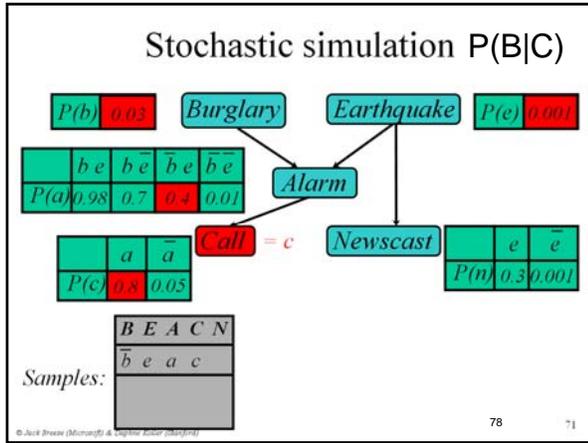
B	E	A	C	N

© Jack Breese (Microsoft) & Daphne Koller (Stanford)

71

71





Stochastic simulation P(B|C)

$P(b) = 0.03$ (Burglary) $P(e) = 0.001$ (Earthquake)
 $P(a|be) = 0.98$, $P(a|b\bar{e}) = 0.7$, $P(a|\bar{b}e) = 0.4$, $P(a|\bar{b}\bar{e}) = 0.01$ (Alarm)
 $P(c|a) = 0.8$, $P(c|\bar{a}) = 0.05$ (Call) $P(n|e) = 0.3$, $P(n|\bar{e}) = 0.001$ (Newscast)

Samples:

B	E	A	C	N
b	e	a	c	n
⋮	⋮	⋮	⋮	⋮

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

© Jack Breese (MIT) & Daphne Koller (Stanford) 84 71

Rejection Sampling

- Sample from the prior
 - reject if do not match the evidence
- Returns consistent posterior estimates
- Hopelessly expensive if P(e) is small
 - P(e) drops exponentially with num of evidence vars

85

Likelihood Weighting

- Idea:
 - fix evidence variables
 - sample only non-evidence variables
 - weight each sample by the likelihood of evidence

86

Likelihood weighting P(B|C)

$P(c|a) = 0.8$, $P(c|\bar{a}) = 0.05$
 $P(\bar{c}|\bar{a}) = 0.2$, $P(\bar{c}|a) = 0.95$ (Call)

Samples:

B	E	A	C	N

© Jack Breese (MIT) & Daphne Koller (Stanford) 87 72

Likelihood weighting P(B|C)

$P(c|a) = 0.8$, $P(c|\bar{a}) = 0.05$
 $P(\bar{c}|\bar{a}) = 0.2$, $P(\bar{c}|a) = 0.95$ (Call)

Samples:

B	E	A	C	N
b				

© Jack Breese (MIT) & Daphne Koller (Stanford) 88 73

Likelihood weighting P(B|C)

$P(c|a) = 0.8$, $P(c|\bar{a}) = 0.05$
 $P(\bar{c}|\bar{a}) = 0.2$, $P(\bar{c}|a) = 0.95$ (Call)

Samples:

B	E	A	C	N
b	e			

© Jack Breese (MIT) & Daphne Koller (Stanford) 89 74

Likelihood weighting P(B|C)

	a	\bar{a}
$P(c)$	0.8	0.05
$P(\bar{c})$	0.2	0.95

Samples:

B	E	A	C	N
b	e	a		

90 72

Likelihood weighting P(B|C)

	a	\bar{a}
$P(c)$	0.8	0.05
$P(\bar{c})$	0.2	0.95

Samples:

B	E	A	C	N
b	e	a	c	

91 72

Likelihood weighting P(B|C)

	a	\bar{a}
$P(c)$	0.8	0.05
$P(\bar{c})$	0.2	0.95

Samples:

B	E	A	C	N
b	e	a	c	\bar{n}

92 72

Likelihood weighting P(B|C)

	a	\bar{a}
$P(c)$	0.8	0.05
$P(\bar{c})$	0.2	0.95

Samples:

B	E	A	C	N	weight
b	e	a	c	\bar{n}	0.8

93 72

Likelihood weighting P(B|C)

	a	\bar{a}
$P(c)$	0.8	0.05
$P(\bar{c})$	0.2	0.95

Samples:

B	E	A	C	N	weight
b	e	a	c	\bar{n}	0.8
b	e	\bar{a}	c	n	0.05

94 72

Likelihood weighting P(B|C)

	a	\bar{a}
$P(c)$	0.8	0.05
$P(\bar{c})$	0.2	0.95

Samples:

B	E	A	C	N	weight
b	e	a	c	\bar{n}	0.8
b	e	\bar{a}	c	n	0.05
					⋮

95 72

Likelihood weighting P(B|C)

	a	\bar{a}
$P(c)$	0.8	0.05
$P(\bar{c})$	0.2	0.95

Samples: **Call** = c

B	E	A	C	N	weight
b	e	a	c	\bar{n}	0.8
b	e	a	c	n	0.05
⋮					

$$P(b|c) = \frac{\text{weight of samples with } B=b}{\text{total weight of samples}}$$

© Jack Breese (Microsoft) & Daphne Koller (Stanford) 96

Likelihood Weighting

- Sampling probability: $S(z, e) = \prod_i P(z_i | \text{Parents}(Z_i))$
 - Neither prior nor posterior
- Wt for a sample $\langle z, e \rangle$: $w(z, e) = \prod_i P(e_i | \text{Parents}(E_i))$
- Weighted Sampling probability $S(z, e)w(z, e)$

$$= \prod_i P(z_i | \text{Parents}(Z_i)) \prod_i P(e_i | \text{Parents}(E_i))$$

$$= P(z, e)$$

→ returns consistent estimates
performance degrades w/ many evidence vars

- a few samples get majority of the weight
- late occurring evidence vars don't guide sample generation

97

MCMC with Gibbs Sampling

- Fix the values of observed variables
- Set the values of all non-observed variables randomly
- Perform a random walk through the space of complete variable assignments. On each move:
 - Pick a variable X
 - Calculate $\Pr(X=\text{true} | \text{all other variables})$
 - Set X to true with that probability
- Repeat many times. Frequency with which any variable Y is true = its posterior probability.
- Converges to true posterior when frequencies stop changing significantly
 - stable distribution, mixing

98

Given Markov Blanket, X is Independent of All Other Nodes

$$MB(X) = \text{Par}(X) \cup \text{Childs}(X) \cup \text{Par}(\text{Childs}(X))$$

© D. Weld and D. Fox 99

Markov Blanket Sampling

- How to calculate $\Pr(X=\text{true} | \text{all other variables})$?
- Recall: a variable is independent of all others given it's Markov Blanket
 - parents
 - children
 - other parents of children
- So problem becomes calculating $\Pr(X=\text{true} | MB(X))$
 - Fortunately, it is easy to solve exactly

$$P(X) = \alpha P(X | \text{Parents}(X)) \prod_{Y \in \text{Childs}(X)} P(Y | \text{Parents}(Y))$$

100

Example

$$P(X) = \alpha P(X | \text{Parents}(X)) \prod_{Y \in \text{Childs}(X)} P(Y | \text{Parents}(Y))$$

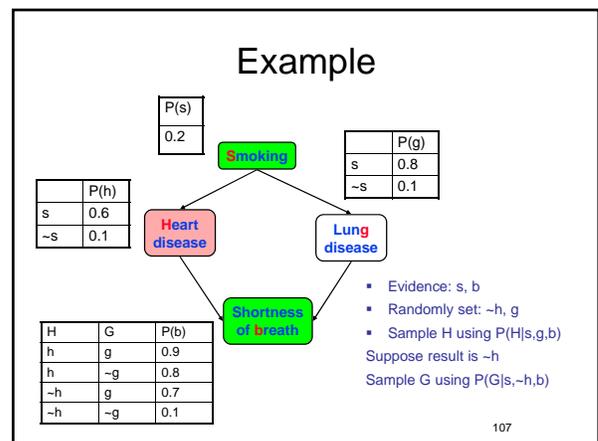
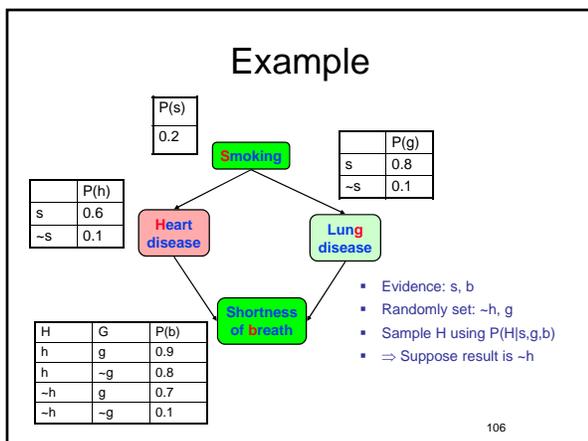
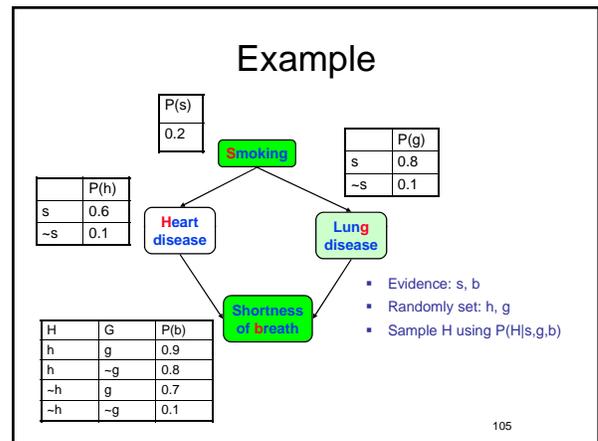
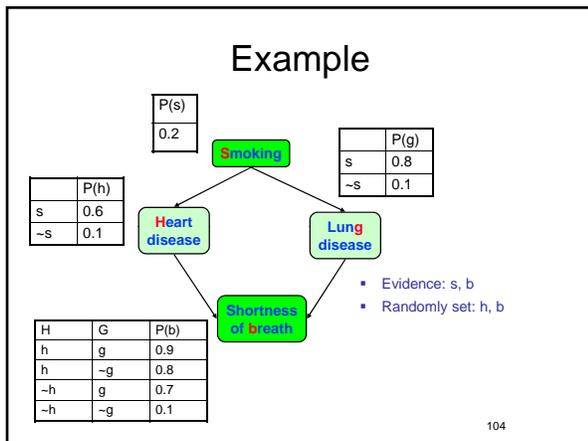
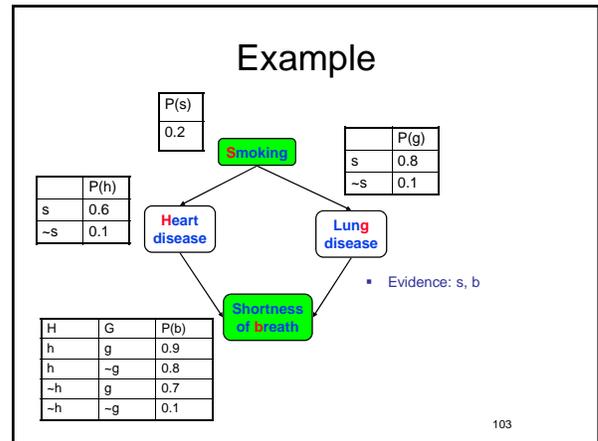
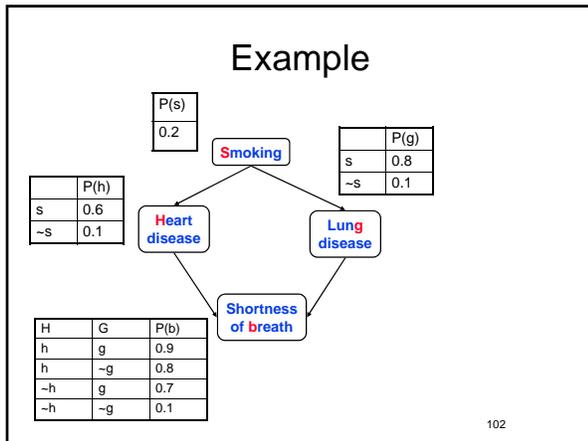
$$P(X | A, B, C) = \frac{P(X, A, B, C)}{P(A, B, C)}$$

$$= \frac{P(A)P(X|A)P(C)P(B|X, C)}{P(A, B, C)}$$

$$= \left[\frac{P(A)P(C)}{P(A, B, C)} \right] P(X|A)P(B|X, C)$$

$$= \alpha P(X|A)P(B|X, C)$$

101



Example

P(s)	
0.2	

P(g)	
s	0.8
~s	0.1

P(h)	
s	0.6
~s	0.1

P(b)		
h	g	0.9
h	~g	0.8
~h	g	0.7
~h	~g	0.1

- Evidence: s, b
- Randomly set: ~h, g
- Sample H using P(H|s,g,b)
- Suppose result is ~h
- Sample G using P(G|s,~h,b)
- ⇒ Suppose result is g

108

Example

P(s)	
0.2	

P(g)	
s	0.8
~s	0.1

P(h)	
s	0.6
~s	0.1

P(b)		
h	g	0.9
h	~g	0.8
~h	g	0.7
~h	~g	0.1

- Evidence: s, b
- Randomly set: ~h, g
- Sample H using P(H|s,g,b)
- Suppose result is ~h
- Sample G using P(G|s,~h,b)
- ⇒ Suppose result is g
- Sample G using P(G|s,~h,b)

Example

P(s)	
0.2	

P(g)	
s	0.8
~s	0.1

P(h)	
s	0.6
~s	0.1

P(b)		
h	g	0.9
h	~g	0.8
~h	g	0.7
~h	~g	0.1

- Evidence: s, b
- Randomly set: ~h, g
- Sample H using P(H|s,g,b)
- Suppose result is ~h
- Sample G using P(G|s,~h,b)
- ⇒ Suppose result is g
- Sample G using P(G|s,~h,b)

Gibbs MCMC Summary

$$P(X|E) = \frac{\text{number of samples with } X=x}{\text{total number of samples}}$$

- **Advantages:**
 - No samples are discarded
 - No problem with samples of low weight
 - Can be implemented very efficiently
 - 10K samples @ second
- **Disadvantages:**
 - Can get stuck if relationship between vars is *deterministic*
 - Many variations devised to make MCMC more robust

111

Other inference methods

- **Exact inference**
 - Junction tree
- **Approximate inference**
 - Belief Propagation
 - Variational Methods

112

Outline

- **Probabilistic models**
 - Bayesian Networks (BNs)
 - Independence in BNs
 - Efficient Inference in BNs
 - **Learning**