# (Conditional) Independence

#### Independence

Two variables are **independent** if:  $\forall x, y P(x, y) = P(x)P(y)$ 

We denote this as  $X \! \perp \! \! \! \! \perp Y$ 

## **Conditional Independence**

#### X is **conditionally independent** of Y given Z

if and only if: 
$$\forall x, y, z : P(x, y|z) = P(x|z)P(y|z)$$

or, equivalently, if and only if  $\forall x, y, z : P(x|z, y) = P(x|z)$ 

#### $X \bot\!\!\!\perp Y | Z$

# **Conditional Independence**

Traffic, Umbrella, Raining

XIIY Z 1 Raining



$$T \perp U^{2}$$

$$T \perp U \mid R$$

$$p(T \mid R, U) = p(T \mid R)$$

# **Conditional Independence**

(Smole datector) Fire, Smoke, Alarm

X TTX S





#### Independence vs. Conditional Independence

Rain Traffic Pedestrian holding umbrella Flood in the house Trip cancelled

. . .



P(Traffic | Rain, Umbrella) = P(Traffic | Rain) Conditional Independent

Conditional distribution / independence allows us to model the probability of a certain event only using relevant factors.

# Bayes Net

# **Bayesian Network Example**

Traffic, Umbrella, Raining

P(t, u, r)

= P(r) P(t | r) P(u | r, t) (always hold by chain rule) = P(r) P(t | r) P(u | r)T  $\perp U | R$ 



# **Bayesian Network (BN)**

- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node
  - Suppose a node as m parents, and suppose each random variable can take d different values

XES

- What is the size of the table?
- The BN models the joint probability as

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

$$\# \text{rms} = \int^{m+\ell} d^{m+\ell}$$

# **Bayesian Network Example**

Fire, Smoke, Alarm







# Recap



А⊥В

#### **Example: Car Insurance**



# **Example: Medical Diagnosis**



Marin Prcela et al. Information Gain of Structured Medical Diagnostic Tests -Integration of Bayesian Networks and Ontologies

# **Causality?**

- When Bayes' nets reflect the true causal patterns:
  - Often simpler (nodes have fewer parents) and easier to think about
- BNs need not be causal
  - Sometimes no causal net exists over the domain (especially if variables are missing)
  - Arrows that reflect correlation, but not necessary causality



#### **Causality?**



# **Independence Given Evidence**

**General question**: Are two variables *X*, *Y* independent of each other conditioned on  $Z = \{Z_1, Z_2, ...\}$ ?

Or: Are X and Y "D-separated" by Z?

#### Algorithm

- 1. Consider just the **ancestral subgraph** consisting of X, Y, Z, and their ancestors.
- 2. Add links between any unlinked pair of nodes that share a common child; now we have the so-called **moral graph**.
- 3. Replace all directed links by undirected links.
- 4. If Z blocks all paths between X and Y in the resulting graph, then Z d-separates X and Y.

#### Example

 $\begin{array}{ll} R \bot B & \text{Yes} \\ R \bot B | T \\ R \bot B | T' \end{array}$ 



#### Example





# Example

- Variables:
  - R: Raining
  - T: Traffic
  - D: Roof drips
  - S: I'm sad
- Questions:
  - $T \! \perp \!\!\! \perp D$
  - $T \perp\!\!\!\perp D | R$  Yes  $T \perp\!\!\!\perp D | R, S$



# **Proof Sketch**

**Statement:** If X and Y and separated by Z in the moral graph, then  $X \perp Y \mid Z$ 



The moral graph gives a way to "factorize" the joint distribution of BN. Each clique in the moral graph is a factor.

$$\underbrace{\mathsf{P}(\mathsf{a}) \; \mathsf{P}(\mathsf{b}) \; \mathsf{P}(\mathsf{c}) \; \mathsf{P}(\mathsf{d} \mid \mathsf{a}, \mathsf{b}, \mathsf{c})}_{\phi(\mathsf{a}, \mathsf{b}, \mathsf{c}, \mathsf{d})} \underbrace{\mathsf{P}(\mathsf{e}) \; \mathsf{P}(\mathsf{f} \mid \mathsf{d}, \mathsf{e})}_{\phi(\mathsf{d}, \mathsf{e}, \mathsf{f})} = \phi(\mathsf{a}, \mathsf{b}, \mathsf{c}, \mathsf{d}) \; \phi(\mathsf{d}, \mathsf{e}, \mathsf{f})$$

#### **Proof Sketch**

**Statement:** If X and Y and separated by Z in the moral graph, then  $X \perp Y \mid Z$ 



#### **Structure Implications**

• Given a Bayes net structure, can run d-separation algorithm to build a complete list of conditional independences that are necessarily true of the form

$$X_i \perp \perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$$

• This list determines the set of probability distributions that can be represented

# **Topology Limits Distributions**

 $\{X \perp\!\!\!\perp Y, X \perp\!\!\!\perp Z, Y \perp\!\!\!\perp Z,$ 

 $X \perp\!\!\!\perp Z \mid Y, X \perp\!\!\!\perp Y \mid Z, Y \perp\!\!\!\perp Z \mid X \}$ 

(Z)

 $\{X \perp\!\!\!\perp Z \mid Y\}$ 

- Given some graph topology G, only certain joint distributions can be encoded
- The graph structure guarantees certain (conditional) independences
- Adding arcs increases the set of distributions, but has several costs

# **Application: Language Modeling**

• Markov Model

- Probabilistic program: Markov model For each position  $i=1,2,\ldots,n$ : Generate word  $X_i\sim p(X_i\mid X_{i-1})$ 



# **Application: Object Tracking**

• Hidden Markov Model

Probabilistic program: hidden Markov model (HMM)
For each time step  $t = 1, \ldots, T$ :
Generate object location  $H_t \sim p(H_t \mid H_{t-1})$ Generate sensor reading  $E_t \sim p(E_t \mid H_t)$ 



Inference: given sensor readings, where is the object?

# **Application: Topic Modeling**

• Latent Dirichlet Allocation

Probabilistic program: latent Dirichlet allocation
Generate a distribution over topics  $\alpha \in \mathbb{R}^K$ For each position  $i = 1, \ldots, L$ :
Generate a topic  $Z_i \sim p(Z_i \mid \alpha)$ Generate a word  $W_i \sim p(W_i \mid Z_i)$ 



Document classification, information retrieval, customer segmentation, ...

Inference: given a text document, what topics is it about?

# **Exact Inference in Bayesian Networks**

## The "Join" Operation in Bayesian Network

The BN defines four factors P(A), P(B|A), P(C|A), P(D|B,C)

Α **Join on B:** Combine all factors that involve B Α P(A), P(B|A), P(C|A), P(D|B,C)С В B,D С P(A), P(B,D | A,C), P(C|A)D **Further join on C:** Combine all factors that involve C P(A), P(B,D | A,C), P(C|A)Α  $P(A), P(B,C,D \mid A)$ B,C,D

#### **Exercise**



BCC

What are the factors after joining on B?

P(A) P(B|A) p(C(A,B) p(b|B,C)) $P(\beta, c, \rho | A)$ 

P(b,a|e) = P(5)P(a|b,e)

#### **Exercise**



# **Review: Inference by Enumeration**

General case:

• Evidence variables: 
$$E_1 \dots E_k = e_1 \dots e_k$$
  
• Query\* variable:  $Q$   
• Hidden variables:  $H_1 \dots H_r$   
 $P(Q|e_1 \dots e_k) = ?$   
 $P(\subseteq I, \dots, E_K, Q, H_1, A, H_Y)$ 

#### **Inference by Enumeration**

Step 1. Select the entries consistent with the evidenceStep 2. Sum out H to get joint probability of Query and evidenceStep 3. Normalize

# **Inference by Enumeration**

Step 0. Create a joint probability table

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

В	Е	А	J	Μ	P(B,E,A,J,M)	
Т	Т	Т	Т	Т	0.001 * 0.002 * 0.95 * 0.90 * 0.70	
Т	Т	Т	Т	F	0.001 * 0.002 * 0.95 * 0.90 * 0.30	
Т	Т	Т	F	Т	0.001 * 0.002 * 0.95 * 0.10 * 0.70	
F	F	F	F	F	0.999 * 0.998 * 0.999 * 0.95 * 0.99	



$$P(B | +j, +m) = ?$$

# **Step 0: Create a Joint Probability Table**

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



Join on A

В

Ε

Α

#### **Inference by Enumeration**

Step 1. Select the entries consistent with the evidence



В	Е	А	J	Μ	P(B,E,A,J,M)	
Т	Т	Т	Т	Т	0.001 * 0.002 * 0.95 * 0.90 * 0.70	
Т	Т	F	Т	Т	0.001 * 0.002 * 0.05 * 0.05 * 0.01	
Т	F	Т	Т	Т	0.001 * 0.998 * 0.94 * 0.90 * 0.70	
Т	F	F	Т	Т	0.001 * 0.998 * 0.06 * 0.05 * 0.01	
F	Т	Т	Т	Т	0.999 * 0.002 * 0.29 * 0.90 * 0.70	
F	Т	F	Т	Т	0.999 * 0.002 * 0.71 * 0.05 * 0.01	
F	F	Т	Т	Т	0.999 * 0.998 * 0.001 * 0.90 * 0.70	
F	F	F	Т	Т	0.999 * 0.998 * 0.999 * 0.05 * 0.01	

$$P(B | +j, +m) = ?$$

#### **Inference by Enumeration**

**Step 2.** Sum out hidden variable to get joint probability of query and evidence (Marginalize)



P( B | +j, +m) = ?

В	J	Μ	P(B,J,M)
Т	Т	Т	0.0006
F	Т	Т	0.0015
### **Inference by Enumeration**

#### Step 3. Normalize



#### **Inference by Enumeration?**



#### How did we do Inference by Enumeration?

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

Α



### How did we do Inference by Enumeration?



### Improving the Algorithm

- First improvement: Instead of eliminating rows inconsistent with the evidence at the end, we will only keep rows consistent with evidence from the beginning.
- Second improvement: Instead of marginalize all hidden variables at the end after joining all variables, we will interleave joining and marginalization.

## **Improving the Algorithm**

 $P(B | +\dot{\sigma}, +m)$ 

#### **Inference by Enumeration**



#### **Query:** P( B | +j, +m) = ?



#### **Query:** P( B | +j, +m) = ?



with the evidence

**Query:** P( B | +j, +m) = ?



**Query:** P( B | +j, +m) = ?



**Query:** P( B | +j, +m) = ?



**Query:** P( B | +j, +m) = ?



We can then get P(B | +j, +m) by normalizing this table

**Query:** P( B | +j, +m) = ?

Can be done in different orders





- Start with initial factors but instantiated by evidence
- While there are still hidden variables:
  - Pick a hidden variable X
  - Join all factors mentioning X
  - Eliminate (sum out) X (i.e., marginalize X)
- Join all the remaining factors
- Normalize

## **Ordering of the Join and Eliminate?**

- The time and space of variable elimination are dominate by the **size of the largest factor** constructed during the algorithm.
- It's hard to determine the optimal ordering
  - Heuristics: Choose the variable that minimize the size of the next factor to be constructed.



# **Approximate Inference in Bayesian Networks**

## Sampling

- Basic idea
  - Draw N samples from a sampling distribution S
  - Compute an approximate posterior probability
  - Show this converges to the true probability P

- Why sample?
  - Often very fast to get a decent approximate answer
  - The algorithms are very simple and general (easy to apply to fancy models)
  - They require very little memory (*O*(*n*))
  - They can be applied to large models, whereas exact algorithms blow up

## **Sampling in Bayes nets**

- Prior sampling
- Rejection sampling
- Likelihood weighting
- Gibbs sampling

### **Prior Sampling**



### **Prior Sampling**

For *i*=1, 2, ..., *n* (in topological order): Sample  $X_i$  from P( $X_i | parents(X_i)$ ) Return ( $x_1, x_2, ..., x_n$ )

## **Prior Sampling**

• This process generates samples with probability:

 $S_{PS}(x_1,\ldots,x_n) = \prod_i P(x_i \mid parents(X_i)) = P(x_1,\ldots,x_n)$ 

...i.e. the BN's joint probability

- Let the number of samples of an event be  $N_{PS}(x_1,...,x_n)$
- Estimate from N samples is  $Q_N(x_1,...,x_n) = N_{PS}(x_1,...,x_n)/N$
- Then  $\lim_{N\to\infty} Q_N(x_1,...,x_n) = \lim_{N\to\infty} N_{PS}(x_1,...,x_n)/N$ =  $S_{PS}(x_1,...,x_n)$ =  $P(x_1,...,x_n)$
- I.e., the sampling procedure is *consistent*

#### Example

- We'll get a bunch of samples from the BN:
  - C, ¬S, r, W C, S, r, W ¬C, S, r, ¬W C, ¬S, r, W ¬C, ¬S, ¬r, W
- If we want to know P(W)
  - We have counts  $\langle w:4, \neg w:1 \rangle$
  - Normalize to get *P(W)* = <w:0.8, ¬w:0.2>
  - This will get closer to the true distribution with more samples



## **Rejection Sampling**

- A simple application of prior sampling for estimating conditional probabilities
  - Let's say we want  $P(C | r, w) = \alpha P(C, r, w)$
  - For these counts, samples with ¬*r* or ¬*w* are not relevant
  - So count the C outcomes for samples with *r*, *w* and reject all other samples
- This is called *rejection sampling* 
  - It is also consistent for conditional probabilities (i.e., correct in the limit)



C, ¬S, r, W \_\_C, \_\_S, ¬r, ¬W \_\_C, \_\_S, \_\_r \_\_C, ¬S, r, W

## **Rejection Sampling**

```
Input: evidence e_1, ..., e_k

For i=1, 2, ..., n

Sample X_i from P(X_i | parents(X_i))

If x_i not consistent with evidence

Reject: Return, and no sample is generated in this cycle

Return (x_1, x_2, ..., x_n)
```

- Problem with rejection sampling:
  - If evidence is unlikely, rejects lots of samples
  - Evidence not exploited as you sample
  - Consider P(*Shape*|*Color=blue*)

- Idea: fix evidence variables, sample the rest
  - Problem: sample distribution not consistent!
  - Solution: *weight* each sample by probability of evidence variables given parents



pyramid, green pyramid, red sphere, blue cube, red sphere, green



pyramid, blue pyramid, blue sphere, blue cube, blue sphere, blue



```
Input: evidence e_1, \dots, e_k
w = 1.0
for i=1, 2, ..., n
   if X_i is an evidence variable
        X_i = observed value<sub>i</sub> for X_i
        Set w = w^* P(x_i | parents(X_i))
   else
        Sample x_i from P(X_i | parents(X_i))
return (x_1, x_2, ..., x_n), w
```

## Likelihood Weighting is Consistent

• Sampling distribution if Z sampled and e fixed evidence

 $S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{j} P(z_j \mid parents(Z_j))$ 

• Now, samples have weights

 $w(\mathbf{z}, \mathbf{e}) = \prod_{k} P(e_{k} | parents(E_{k}))$ 

• Together, weighted sampling distribution is consistent

 $S_{WS}(\mathbf{z}, \mathbf{e}) \cdot w(\mathbf{z}, \mathbf{e}) = \prod_{j} P(z_{j} \mid parents(Z_{j})) \prod_{k} P(e_{k} \mid parents(E_{k}))$  $= P(\mathbf{z}, \mathbf{e})$ 



- Likelihood weighting is good
  - All samples are used
  - The values of *downstream* variables are influenced by *upstream* evidence



- Likelihood weighting still has weaknesses
  - The values of *upstream* variables are unaffected by *downstream* evidence
    - E.g., suppose evidence is a video of a traffic accident
  - With evidence in k leaf nodes, weights will be  $O(2^{-k})$
  - With high probability, one lucky sample will have much larger weight than the others, dominating the result
- We would like each variable to "see" *all* the evidence!

-b,-e,-a,+j:,-m~ -b .- l .- a , - ; ,- m ~ b = e , a = - s = m b = e , a + = + m 



**Query:** P( B | +j, +m) = ?

## **Gibbs sampling**

- Ideas
  - States are complete assignments to all variables
    - like *local search* (for constraint satisfaction problems)
  - Evidence variables remain fixed, other variables change
  - To generate the next state, pick a variable and sample a value for it conditioned on all the other variables: X<sub>i</sub>' ~ P(X<sub>i</sub> | x<sub>1</sub>,..,x<sub>i-1</sub>,x<sub>i+1</sub>,..,x<sub>n</sub>)
    - Will tend to move towards states of higher probability, but can go down too
    - In a Bayes net,  $P(X_i | x_1, ..., x_{i+1}, ..., x_n) = P(X_i | markov_blanket(X_i))$
- Theorem: Gibbs sampling is consistent\*

## Gibbs Sampling Example: P( S | +r)

- Step 1: Fix evidence
  - R = +r



- Step 2: Initialize other variables
  - Randomly



- Steps 3: Repeat
  - Choose a non-evidence variable X
  - Resample X from P(X | all other variables)



#### **Advantages of MCMC**



- Samples soon begin to reflect all the evidence in the network
- Eventually they are being drawn from the true posterior!

#### Car Insurance: *P*(*PropertyCost* | *e*)



#### Car Insurance: *P*(*PropertyCost* | *e*)


## Gibbs sampling algorithm

• Repeat many times: Sample a non-evidence variable  $X_i$  from

l

involves X;

- $P(X_i | x_1, ..., x_{i-1}, x_{i+1}, ..., x_n)$
- $= P(X_i | \text{Markov\_blanket}(X_i))$
- =  $\alpha P(X_i | \text{Parents}(X_i)) \prod_i P(y_i | \text{Parents}(Y_i))$
- Markov\_blanket(X<sub>i</sub>) includes
  - $X_i$ 's parents
  - $X_i$ 's children
  - $X_i$ 's children's parent



## **Efficient Resampling of One Variable**

$$P(S | +c, +r, -w)$$

$$= \frac{P(S, +c, +r, -w)}{P(+c, +r, -w)} (S = +s, -s)$$

$$= \frac{P(S, +c, +r, -w)}{\sum P(s, +c, +r, -w)}$$

$$= \frac{P(S, +c, +r, -w)}{\sum P(s, +c, +r, -w)}$$

$$= \frac{P(FC) P(S | +c) P(-w | s, -c, +r)}{\sum P(FC) P(s | +c) P(-w | s, -c, +r)}$$



## **Gibbs Sampling in practice**

- The most commonly used method for large Bayes nets
  - See, e.g., BUGS, JAGS, STAN, infer.net, BLOG, etc.
- Can be *compiled* to run very fast
  - Eliminate all data structure references, just multiply and sample
  - ~100 million samples per second on a laptop
- Can run asynchronously in parallel (one processor per variable)

## **Bayes' Net Inference Summary**

- Exact Inference
  - Inference by Enumeration
  - Variable Elimination
- Approximate Inference
  - Prior sampling
  - Rejection sampling
  - Likelihood weighting
  - Gibbs sampling