

# Introduction to Artificial Intelligence

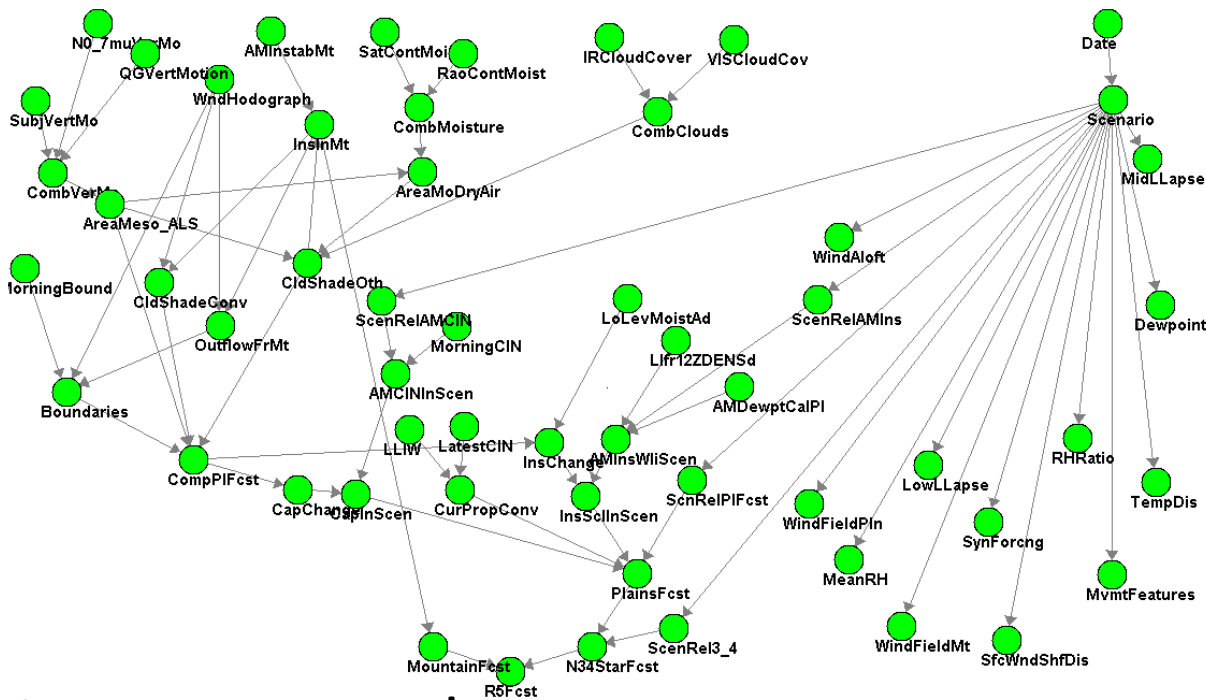
## Lecture 13 – Approximate Inference

CS/CNS/EE 154

Andreas Krause

# Bayesian networks

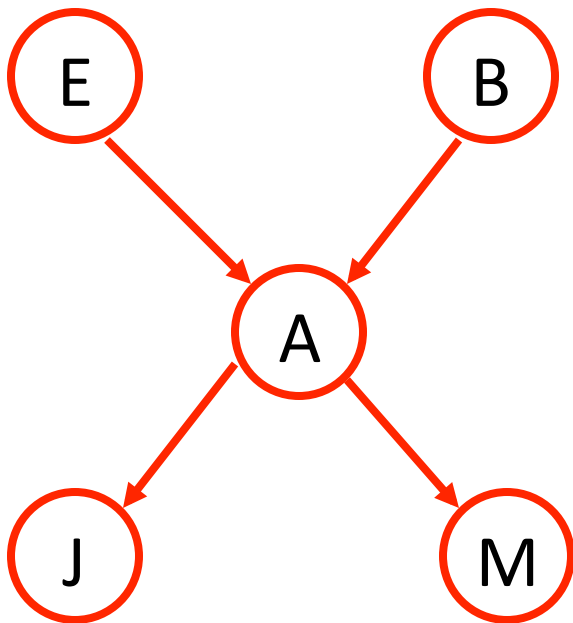
- Compact representation of distributions over large number of variables
- (Often) allows efficient exact inference (computing marginals, etc.)



JavaBayes applet

**HailFinder**  
56 vars  
~ 3 states each  
→ ~ $10^{26}$  terms  
> 10.000 years  
on Top  
supercomputers

# Typical queries: Conditional distribution



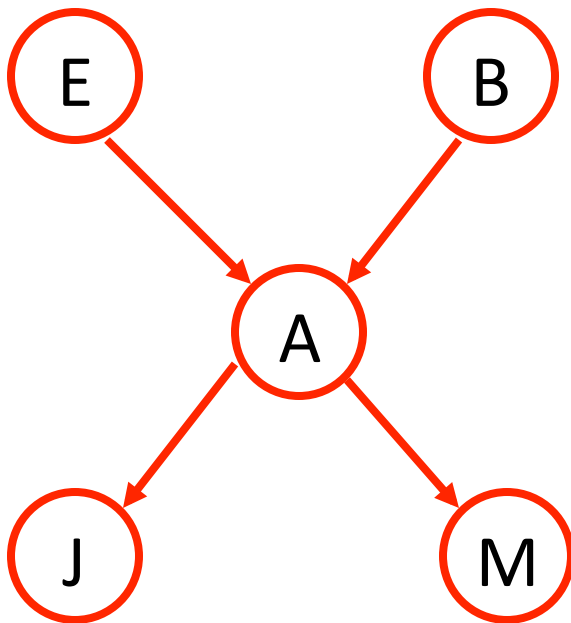
- Compute distribution of some variables given values for others

$$P(E \mid J=T) \quad ?$$

$$P(E, B \mid J=T, M=F) \quad ?$$

$$= \frac{1}{2} P(E, B, J=T, M=F)$$

# Typical queries: Maximization



- MPE (Most probable explanation):  
Given values for some vars,  
compute most likely assignment to  
all remaining vars

$$(a^*, e^*, b^*) = \underset{e, b, a}{\operatorname{argmax}} P(E=e, B=b, A=a \mid J=T, M=F)$$

- MAP (Maximum a posteriori):  
Compute most likely assignment to  
some variables

More general  
than MPE

$$(e^*, b^*) = \underset{e, b}{\operatorname{argmax}} P(e, b \mid J=T, M=F)$$

# Hardness of inference for general BNs

- Computing conditional distributions:
  - Exact solution: #P-complete
  - NP-hard to obtain any nontrivial approximation  
Eg. NP-hard to obtain  $\hat{P}$  s.t.  $|P - \hat{P}| < \frac{1}{2}$
- Maximization:
  - MPE: NP-complete
  - MAP: NP<sup>PP</sup>-complete
- Inference in general BNs is really hard ☹️

# Inference

- Can exploit structure (conditional independence) to efficiently perform **exact inference** in many practical situations
- For BNs where exact inference is not possible, can use algorithms for **approximate inference** (later)

# Variable elimination algorithm

- Given BN and Query  $P(X \mid \mathbf{E}=\mathbf{e})$
- Choose an ordering of  $X_1, \dots, X_n$
- Set up initial factors:  $f_i = P(X_i \mid \mathbf{Pa}_i)$
- For  $i = 1:n$ ,  $X_i \notin \{X, \mathbf{E}\}$ 
  - Collect all factors  $f$  that include  $X_i$
  - Generate new factor by marginalizing out  $X_i$

$$g = \sum_{x_i} \prod_j f_j$$

- Add  $g$  to set of factors
- Renormalize  $P(x, \mathbf{e})$  to get  $P(x \mid \mathbf{e})$

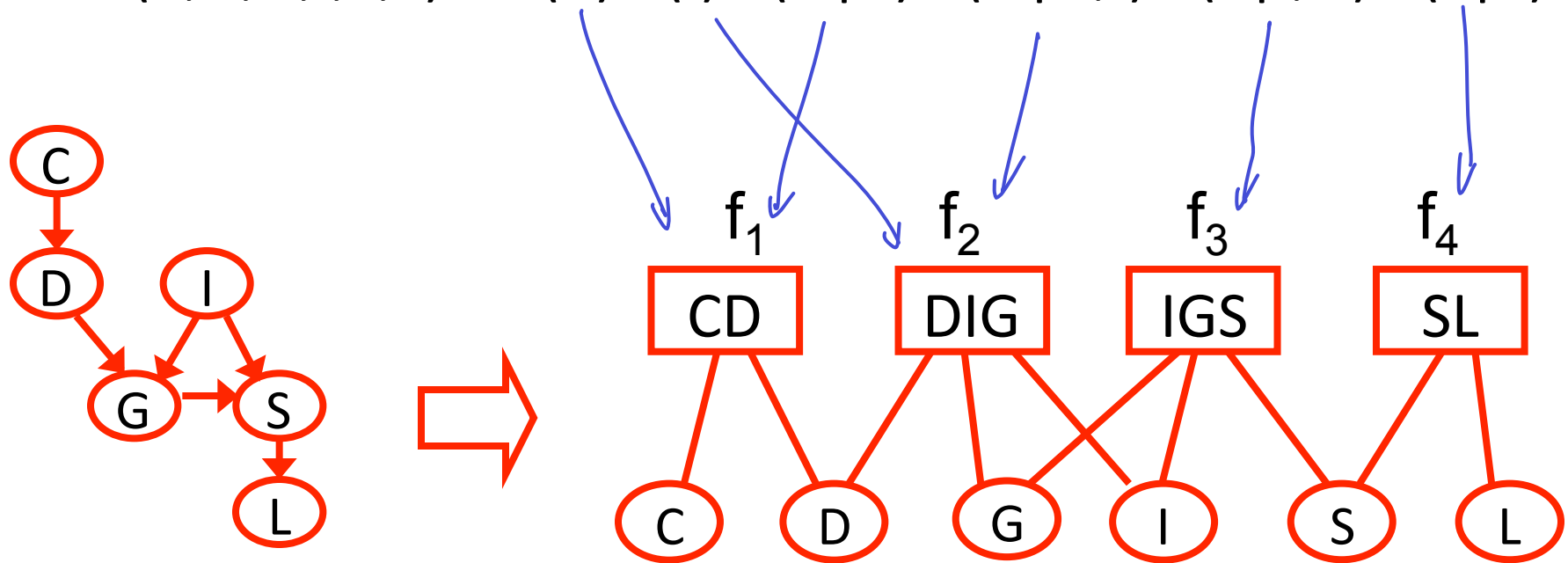
# Reusing computation

- Often, want to compute conditional distributions of many variables, for fixed observations
  - E.g., probability of *Pits* at different locations given observed *Breezes*
  - Repeatedly performing variable elimination is wasteful (many factors are recomputed)
  - Need right data-structure to avoid recomputation
- ➔ Message passing on factor graphs



# Factor graphs

- $P(C,D,G,I,S,L) = P(C) P(I) P(D|C) P(G|D,I) P(S|I,G) P(L|S)$



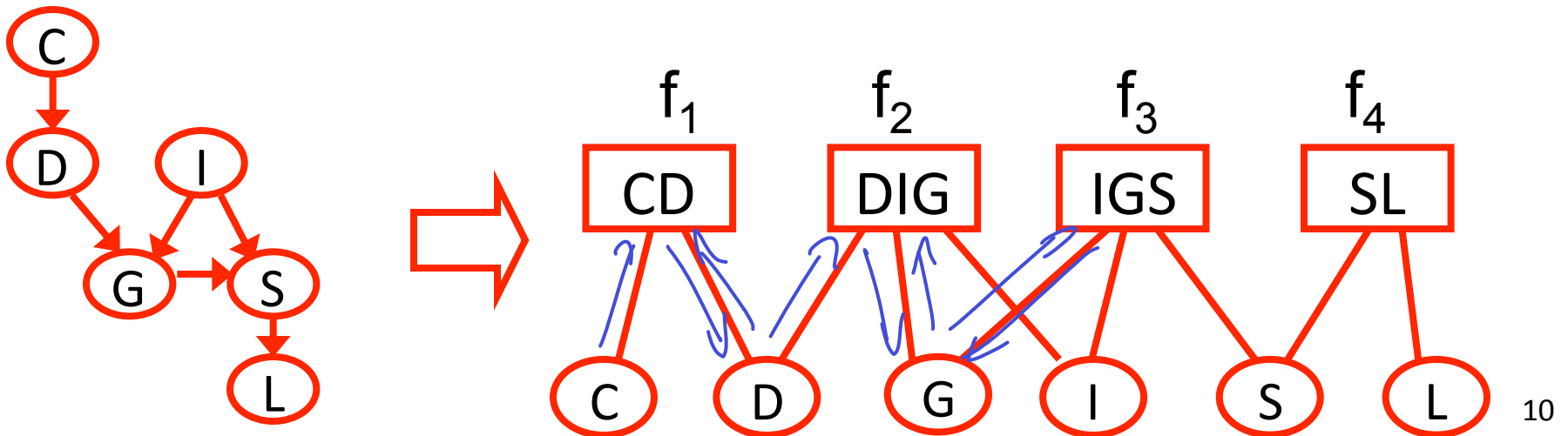
$$f_1(c, d) = P(c) P(d|c)$$

$$f_2(d, i, g) = P(i) P(g|d, i)$$

$$f_3(i, g, s) = P(s|i, g)$$

# Factor graph

- A **factor graph** for a Bayesian network is a bipartite graph consisting of
  - **Variables** and
  - **Factors**
- Each factor is associated with a subset of variables, and all CPDs of the Bayesian network have to be assigned to one of the factor nodes



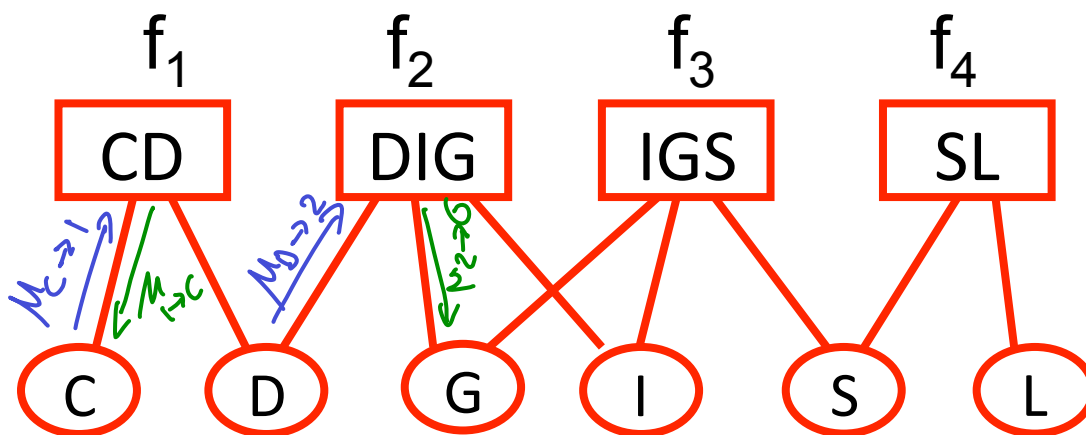
# Sum-product message passing on factor graphs

- Messages from node  $v$  to factor  $u$

$$M_{v \rightarrow u}(x_v) = \prod_{u' \in \mathcal{N}(v) \setminus \{u\}} M_{u' \rightarrow v}(x_v)$$

- Messages from factor  $u$  to node  $v$

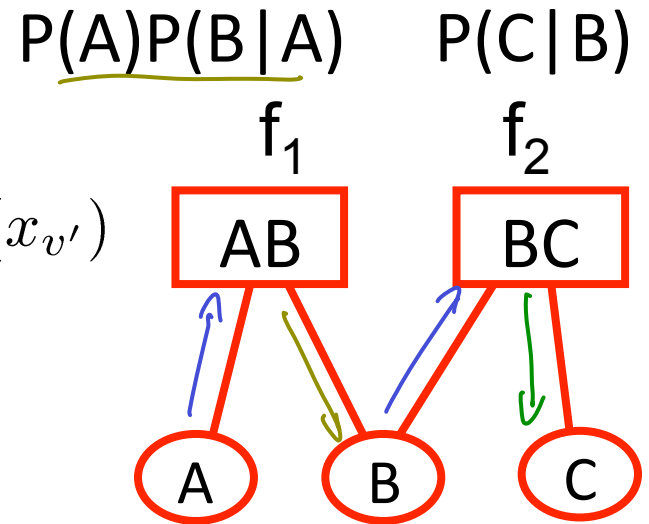
$$M_{u \rightarrow v}(x_v) = \sum_{x_u \sim x_v} f_u(x_u) \prod_{v' \in \mathcal{N}(u) \setminus \{v\}} M_{v' \rightarrow u}(x_{v'})$$



# Example messages

$$\underline{\mu_{v \rightarrow u}(x_v)} = \prod_{u' \in N(v) \setminus \{u\}} \mu_{u' \rightarrow v}(x_v)$$

$$\underline{\mu_{u \rightarrow v}(x_v)} = \sum_{\mathbf{x}_u \sim x_v} f_u(\mathbf{x}_u) \prod_{v' \in N(u) \setminus \{v\}} \mu_{v' \rightarrow u}(x_{v'})$$



$$\mu_{A \rightarrow 1}(a) = 1$$

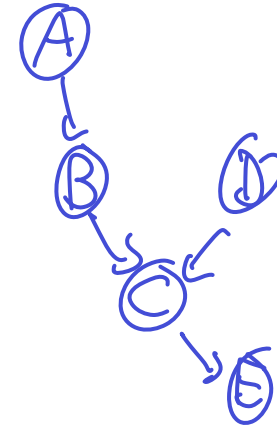
$$\mu_{1 \rightarrow B}(b) = \sum_a f_1(a, b) \cdot \mu_{A \rightarrow 1}(a) = \sum_a \underbrace{P(a)P(b|a)}_{P(a, b)} \cdot 1 = P(b)$$

$$\mu_{B \rightarrow 2}(b) = \mu_{1 \rightarrow B}(b) = P(b)$$

$$\mu_{2 \rightarrow C}(c) = \sum_b f_2(b, c) \cdot \mu_{B \rightarrow 2}(b) = \sum_b \underbrace{P(c|b) \cdot P(b)}_{P(b, c)} = P(c)$$

# Belief propagation on polytrees

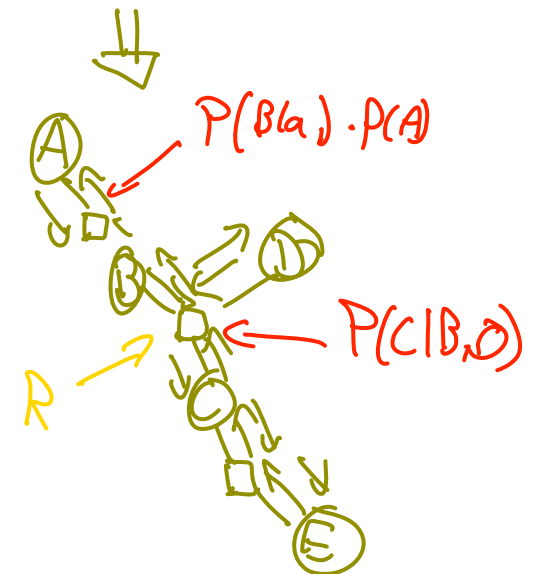
- Belief propagation (aka sum-product) is *exact* for polytree Bayesian networks
  - Factor graph of polytree is a tree
  - Choose one node as root
  - Send messages from leaves to root, and from root to leaves



- After convergence:

$$P(X_v = x_v) = \prod_{u \in N(v)} \mu_{u \rightarrow v}(x_v)$$

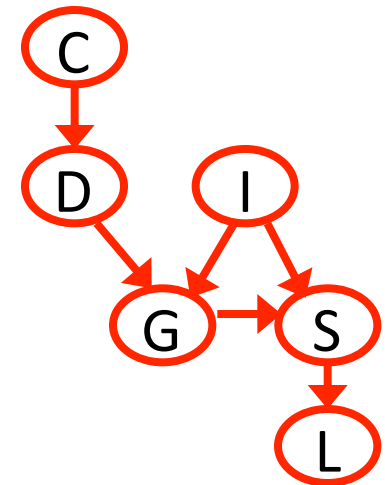
$$P(\mathbf{X}_u = \mathbf{x}_u) = f_u(\mathbf{x}_u) \prod_{v \in N(u)} \mu_{v \rightarrow u}(x_v)$$



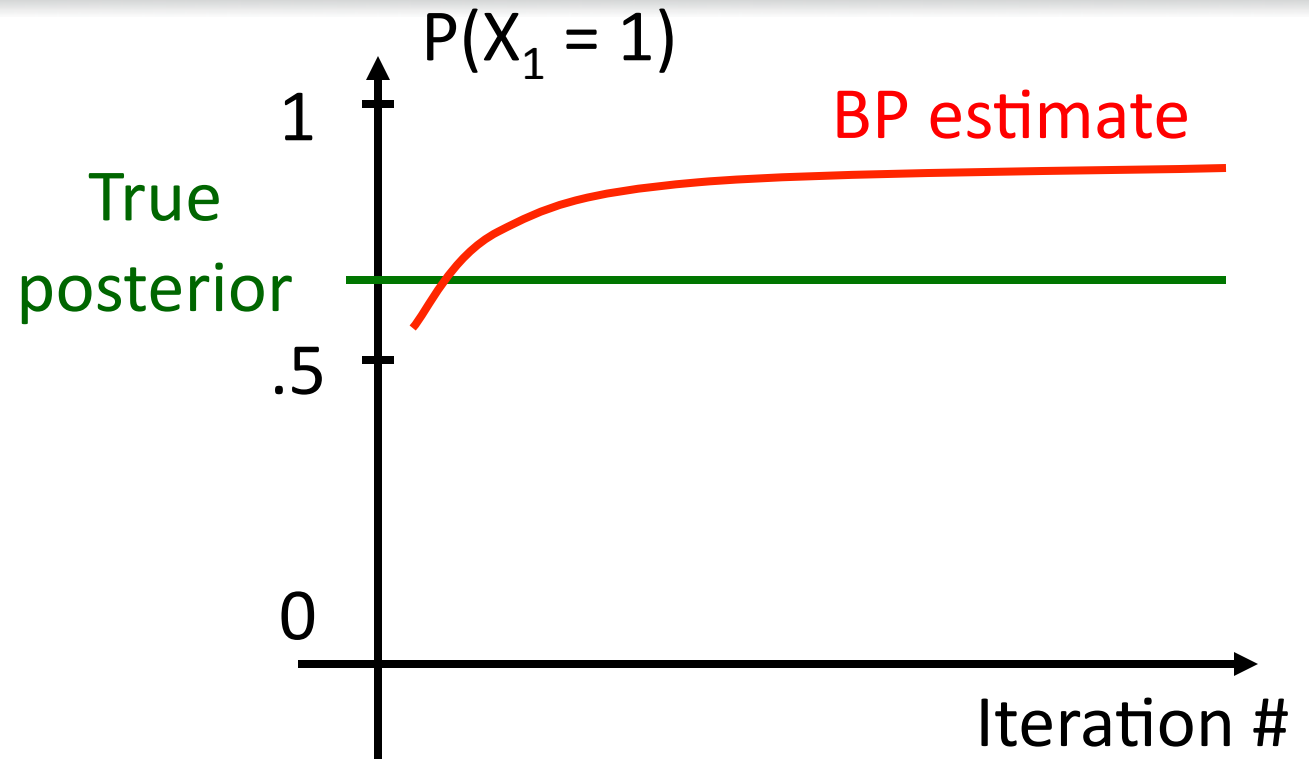
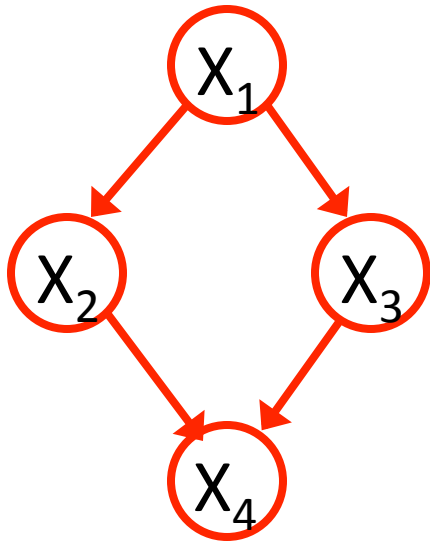
- Thus: immediately have correct values for **all** marginals!<sup>13</sup>

# What if we have loops?

- Can still apply belief propagation even if we have loops
  - Just run it, close your eyes and hope for the best!
  - Use approximation: 
$$P(X_v = x_v) \approx \prod_{u \in N(v)} \mu_{u \rightarrow v}(x_v)$$
- In general, will not converge...
- Even if it converges, may converge to incorrect marginals...
- However, in practice often still useful!
  - E.g., turbo-codes, etc.
- “Loopy belief propagation”



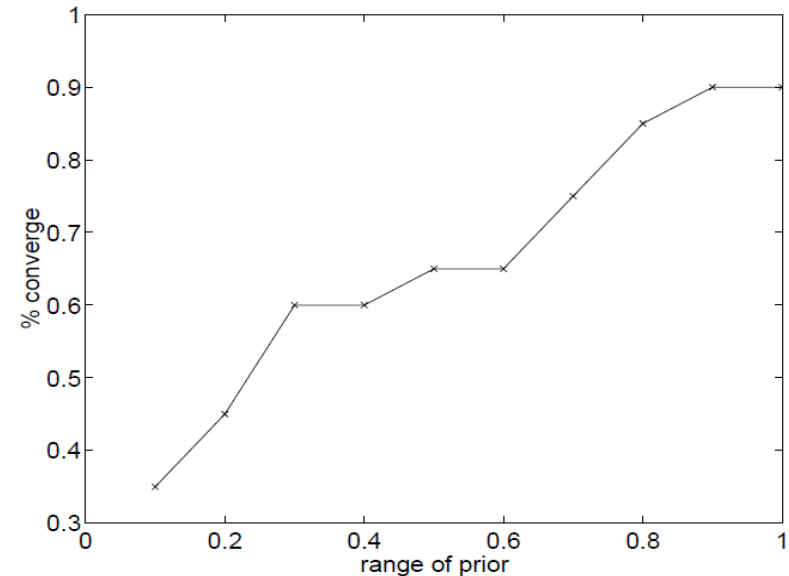
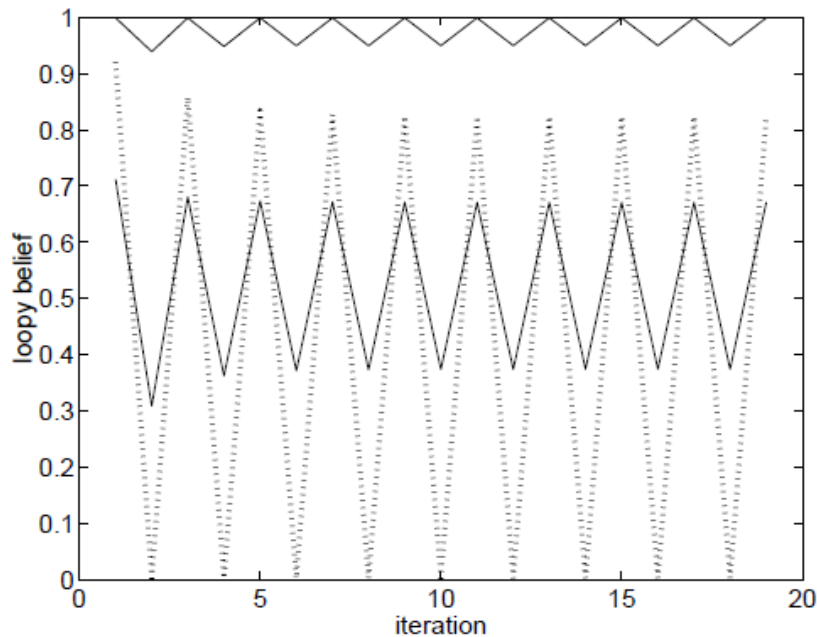
# Behavior of Loopy BP



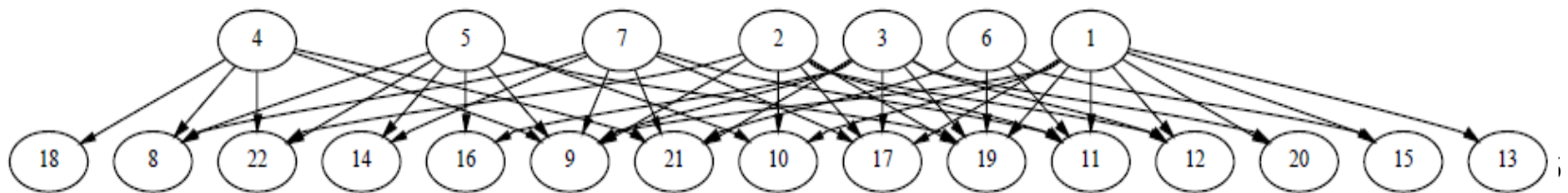
- Loopy BP multiplies same factors multiple times  
→ **BP often overconfident**

# Does Loopy BP always converge?

- No! Can oscillate!
- Typically, oscillation the more severe the more “deterministic” the potentials



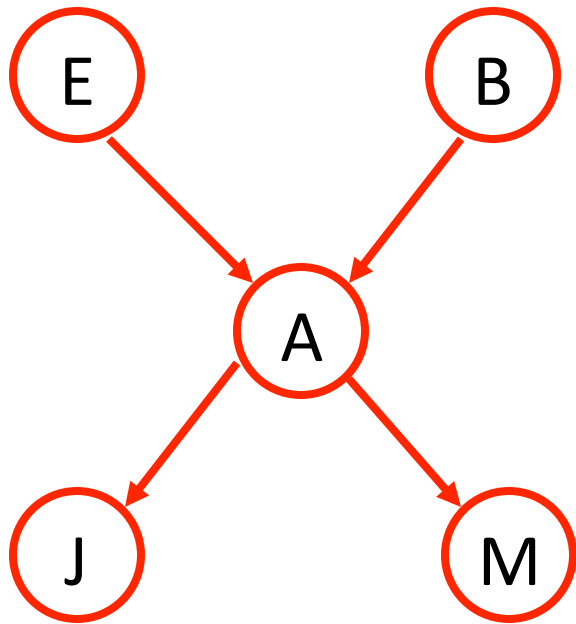
Graphs from K. Murphy UAI '99





# What about MPE queries?

- E.g.,: What's the most likely assignment to the unobserved variables, given the observed ones?



$$(e^*, b^*, a^*) = \underset{e, b, a}{\operatorname{argmax}} P(e, b, a | j, m)$$

$$= \underset{e, b, a}{\operatorname{argmax}} P(e, b, a, j, m)$$

$$\max_{e, b, a} P(e, b, a, j, m)$$

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

$$= \max_a P(j | a) P(m | a) \max_e P(e) \underbrace{\max_b P(b) P(a | e, b)}_{g_b(e, a)}$$

- Use max-product (same as sum-product/BP, but with max instead of sums!)

# Max-product message passing on factor graphs

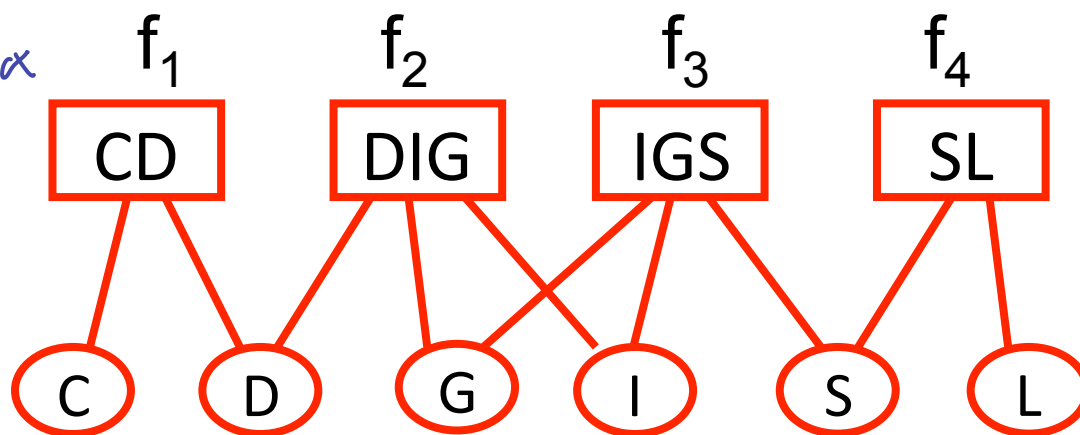
- Messages from nodes to factors

$$\mu_{v \rightarrow u}(x_v) = \prod_{u' \in N(v) \setminus \{u\}} \mu_{u' \rightarrow v}(x_v)$$

- Messages from factors to nodes

$$\mu_{u \rightarrow v}(x_v) = \max_{\mathbf{x}_u \sim x_v} f_u(\mathbf{x}_u) \prod_{v' \in N(u) \setminus \{v\}} \mu_{v' \rightarrow u}(x_{v'})$$

Replace  $\sum$  by  $\max$



# Sampling based inference

- So far: deterministic inference techniques
  - Variable elimination
  - (Loopy) belief propagation
- Will now introduce stochastic approximations
  - Algorithms that “randomize” to compute expectations
  - In contrast to the deterministic methods, guaranteed to converge to right answer (if wait looong enough..)
  - More exact, but slower than deterministic variants

# Computing expectations

- Often, we're not necessarily interested in computing marginal distributions, but certain expectations:
- Moments (mean, variance, ...)

$$\mathbb{E}_P[X^k] = \int x^k P(x) dx$$

- Event probabilities

$$P(X > c) = \mathbb{E}_P[I_{X>c}] = \int [x > c] P(x) dx$$

(in general:  $\mathbb{E}_P(f(x)) = \int P(x) f(x) dx$  for continuous  
=  $\sum_x P(x) f(x)$  for discrete

# Sample approximations of expectations

- $x_1, \dots, x_N$  samples from RV  $X$
- Law of large numbers:

$$\mathbb{E}_P[f(X)] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i)$$

- Hereby, the convergence is with probability 1 (almost sure convergence)
- Finite samples:

$$\mathbb{E}_P(f(x)) \approx \frac{1}{N} \sum_{i=1}^N f(x_i)$$

# How many samples do we need?

- Hoeffding inequality

Suppose  $f$  is bounded in  $[0, C]$ . Then

$$P\left(\left|\mathbb{E}_P[f(X)] - \underbrace{\frac{1}{N} \sum_{i=1}^N f(x_i)}_{\text{Sample avg}}\right| > \varepsilon\right) \leq 2 \exp(-2N\varepsilon^2/C^2)$$

- Thus, probability of error decreases exponentially in  $N$ !

Suppose I want error at most  $\varepsilon$  with prob. at least  $1-\delta$

$$\exp(-2N\varepsilon^2/C^2) \leq \delta/2$$

$$\Leftrightarrow -2N\varepsilon^2/C^2 \leq \ln \delta/2$$

$$\Leftrightarrow 2N\varepsilon^2/C^2 \geq \ln 2/\delta$$

$$\Leftrightarrow N \geq \frac{1}{2} \frac{C^2}{\varepsilon^2} \ln \frac{2}{\delta}$$

- Need to be able to draw samples from  $P$

# Sampling from a Bernoulli distribution

- Most random number generators produce (approximately) uniformly distributed random numbers
- How can we draw samples from  $X \sim \text{Bernoulli}(p)$ ?

$$Y \sim u([0,1])$$

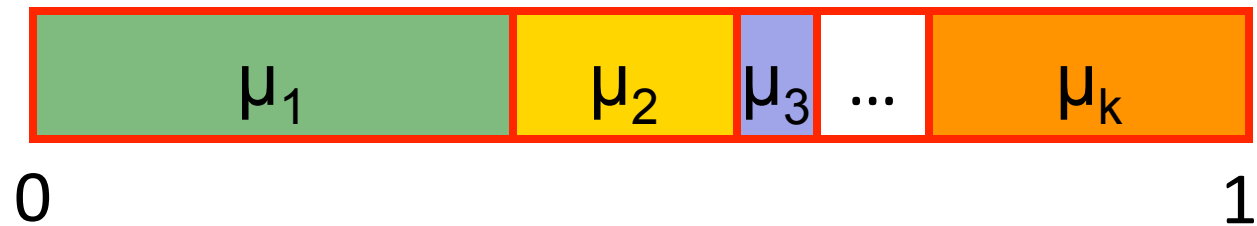
$$X = 1 \quad \text{if } Y < p$$

$$X = 0 \quad \text{if } Y \geq p$$



# Sampling from a Multinomial

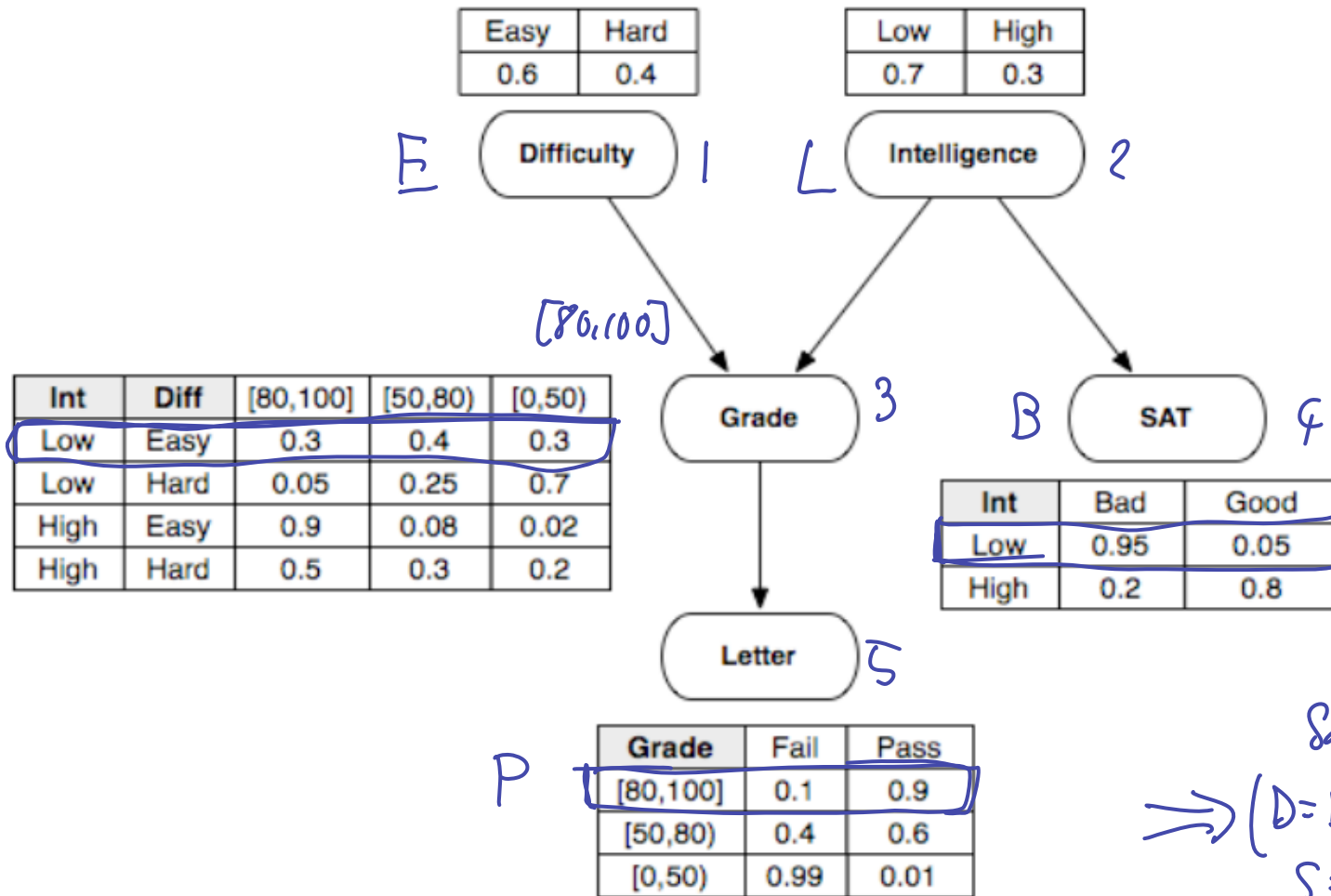
- $X \sim \text{Mult}([\mu_1, \dots, \mu_k])$   
where  $\mu_i = P(X=i)$ ;  $\sum_i \mu_i = 1$



- Function  $g: [0, 1] \rightarrow \{1, \dots, k\}$  assigns state  $g(x)$  to each  $x$
- Draw sample from uniform distribution on  $[0, 1]$
- Return  $g^{-1}(x)$



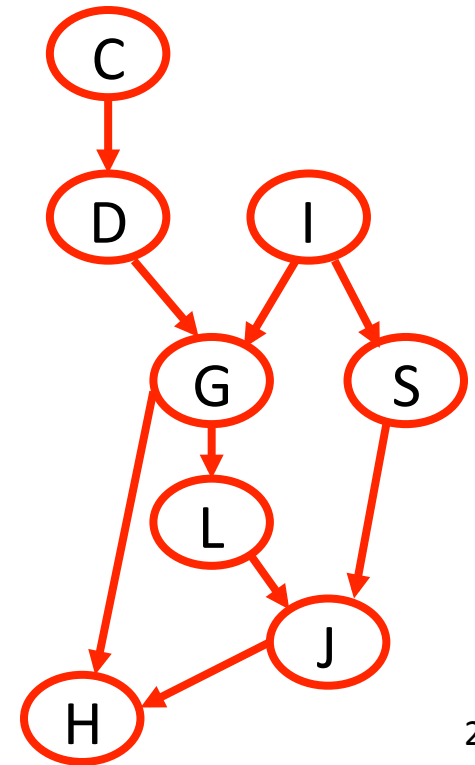
# Forward sampling from a BN



# Monte Carlo sampling from a BN

- Sort variables in topological ordering  $X_1, \dots, X_n$
- For  $i = 1$  to  $n$  do
  - Sample  $x_i \sim P(X_i \mid X_1=x_1, \dots, X_{i-1}=x_{i-1})$

- Works even with loopy models!



# Computing probabilities through sampling

- Want to estimate probabilities
- Draw  $N$  samples from BN

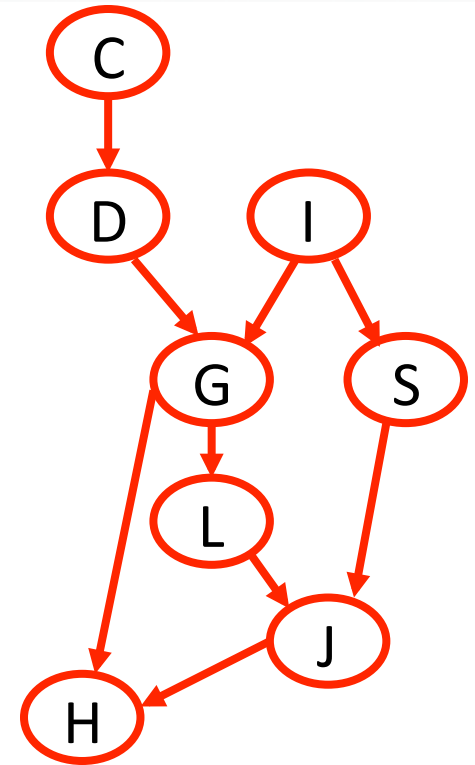
$x_1 \dots x_N$

- Marginals

$$P(L=T) \approx \frac{1}{N} \sum_{i=1}^N \underbrace{[L(x_i)=T]}_{\substack{=1 \text{ if } L(x_i)=T \\ 0 \text{ otherwise}}} = \frac{\text{Count}(L=T)}{N}$$

- Conditionals

$$P(L=T | H=F) = \frac{P(L=T, H=F)}{P(H=F)} \approx \frac{\text{Count}(L=T, H=F)}{\text{Count}(H=F)}$$



# Rejection sampling

- Collect samples over all variables

$$\hat{P}(\mathbf{X}_A = \mathbf{x}_A \mid \mathbf{X}_B = \mathbf{x}_B) \approx \frac{\text{Count}(\mathbf{x}_A, \mathbf{x}_B)}{\text{Count}(\mathbf{x}_B)}$$

- Throw away samples that disagree with  $\mathbf{x}_B$
- Can be problematic if  $P(\mathbf{X}_B = \mathbf{x}_B)$  is rare event

# Sample complexity for probability estimates

- Absolute error:

$$\text{Prob}\left(|\hat{P}(\mathbf{x}) - P(\mathbf{x})| > \varepsilon\right) \leq$$

- Relative error:

$$\text{Prob}\left(\hat{P}(\mathbf{x}) < (1 + \varepsilon)P(\mathbf{x})\right) \leq 2 \exp(-NP(\mathbf{x})\varepsilon^2/3)$$

*if  $P(\mathbf{x})$  exponentially small, need  $N$  exponentially large*

# Sampling from rare events

- Estimating conditional probabilities  $P(X_A | \mathbf{X}_B = \mathbf{x}_B)$  using rejection sampling is hard!
  - The more observations, the unlikelier  $P(\mathbf{X}_B = \mathbf{x}_B)$  becomes
- Want to directly sample from posterior distribution!

# Gibbs sampling

- Start with initial assignment  $\mathbf{x}^{(0)}$  to all variables
- For  $t = 1$  to  $\infty$  do
  - Set  $\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)}$
  - For each variable  $X_i$ 
    - Set  $\mathbf{v}_i =$  values of all  $\mathbf{x}^{(t)}$  except  $x_i$
    - Sample  $x_i^{(t)}$  from  $P(X_i | \mathbf{v}_i)$
- For large enough  $t$ , sampling distribution will be “close” to true posterior distribution!
- **Key challenge:** Computing conditional distributions  $P(X_i | \mathbf{v}_i)$