

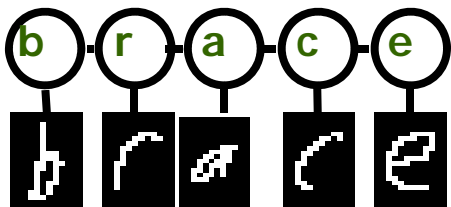
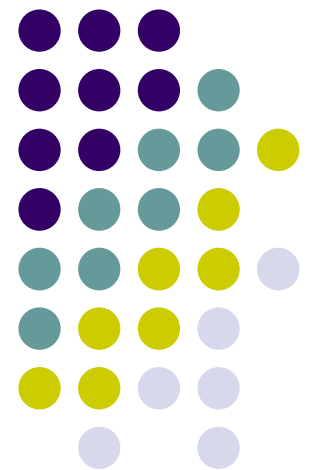
# Machine Learning

10-701, Fall 2016

## Inference and Learning in GM

Eric Xing

Lecture 18, November 9, 2016



Reading: Chap. 8, C.B book



# Inference and Learning

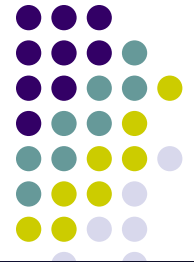
- We now have compact representations of probability distributions: **BN**
- A BN  $M$  describes a unique probability distribution  $P$
- Typical tasks:
  - Task 1: How do we answer **queries** about  $P$ ?
    - We use **inference** as a name for the process of computing answers to such queries
  - Task 2: How do we estimate a **plausible model**  $M$  from data  $D$ ?
    - i. We use **learning** as a name for the process of obtaining point estimate of  $M$ .
    - ii. But for *Bayesian*, they seek  $p(M | D)$ , which is actually an **inference** problem.
    - iii. When not all variables are observable, even computing point estimate of  $M$  need to do **inference** to impute the *missing data*.



# Approaches to inference

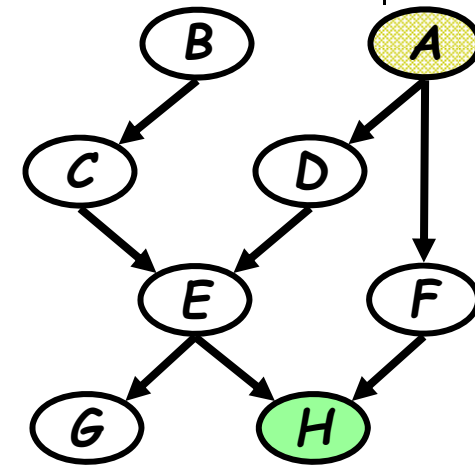
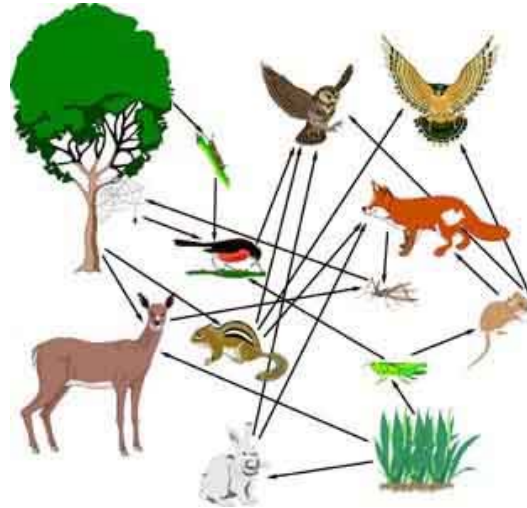
---

- Exact inference algorithms
  - The elimination algorithm
  - Belief propagation
  - The junction tree algorithms (but will not cover in detail here)
- Approximate inference techniques
  - Variational algorithms
  - Stochastic simulation / sampling methods
  - Markov chain Monte Carlo methods



# Marginalization and Elimination

- A food web:



What is the probability that hawks are leaving given that the grass condition is poor?

Query:  $P(h)$

$$P(h) = \sum_g \sum_f \sum_e \sum_d \sum_c \sum_b \sum_a P(a, b, c, d, e, f, g, h)$$



a naïve summation needs to enumerate over an exponential number of terms

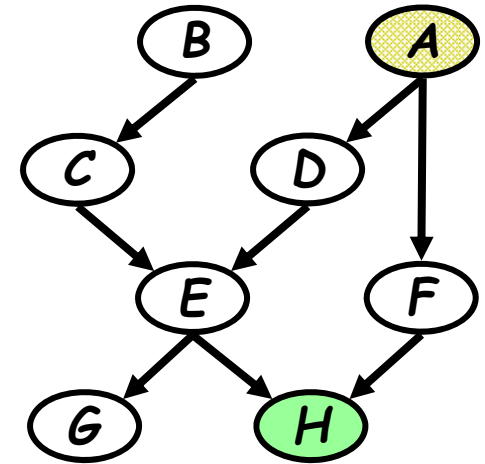
- By chain decomposition, we get

$$= \sum_g \sum_f \sum_e \sum_d \sum_c \sum_b \sum_a P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f)$$

# Variable Elimination



- Query:  $P(A | h)$ 
  - Need to eliminate:  $B, C, D, E, F, G, H$
- Initial factors:
 
$$P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f)$$
- Choose an elimination order:  $H, G, F, E, D, C, B$

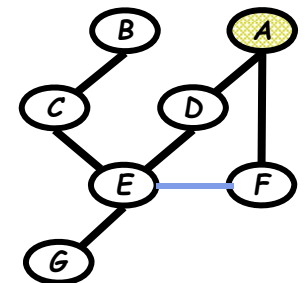


- Step 1:
  - **Conditioning** (fix the evidence node (i.e.,  $h$ ) on its observed value (i.e.,  $\tilde{h}$ )):

$$m_h(e, f) = p(h = \tilde{h} | e, f)$$

- This step is isomorphic to a marginalization step:

$$m_h(e, f) = \sum_h p(h | e, f) \delta(h = \tilde{h})$$

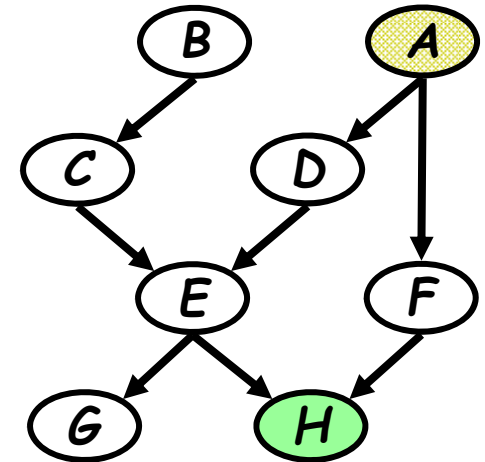




# Example: Variable Elimination

- Query:  $P(B | h)$ 
  - Need to eliminate:  $B, C, D, E, F, G$
- Initial factors:

$$P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f)$$
$$\Rightarrow P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)\underline{m_h(e, f)}$$

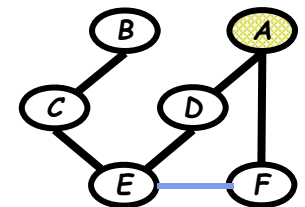


- Step 2: Eliminate  $G$

- compute

$$m_g(e) = \sum_g p(g | e) = 1$$

$$\Rightarrow P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)m_g(e)m_h(e, f)$$
$$= P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)\underline{m_h(e, f)}$$





# Example: Variable Elimination

- Query:  $P(B | h)$ 
  - Need to eliminate:  $B, C, D, E, F$
- Initial factors:

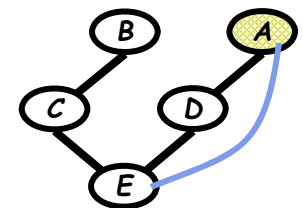
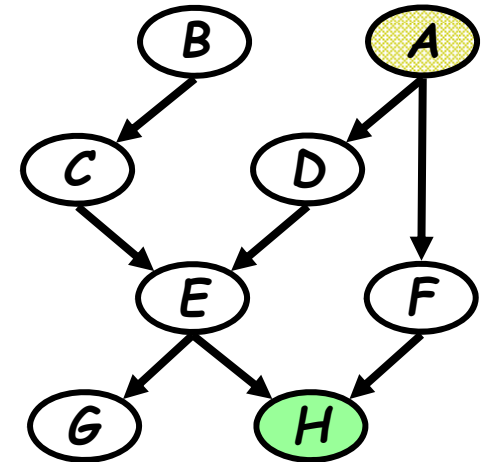
$$\begin{aligned} &P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f) \\ \Rightarrow &P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)m_h(e, f) \\ \Rightarrow &P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)\underline{m_h(e, f)} \end{aligned}$$

- Step 3: Eliminate  $F$

- compute

$$m_f(e, a) = \sum_f p(f | a)m_h(e, f)$$

$$\Rightarrow P(a)P(b)P(c | b)P(d | a)P(e | c, d)\underline{m_f(a, e)}$$





# Example: Variable Elimination

- Query:  $P(B | h)$ 
  - Need to eliminate:  $B, C, D, E$

- Initial factors:

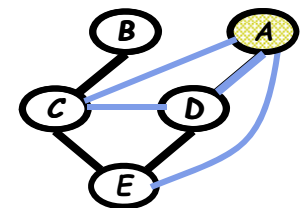
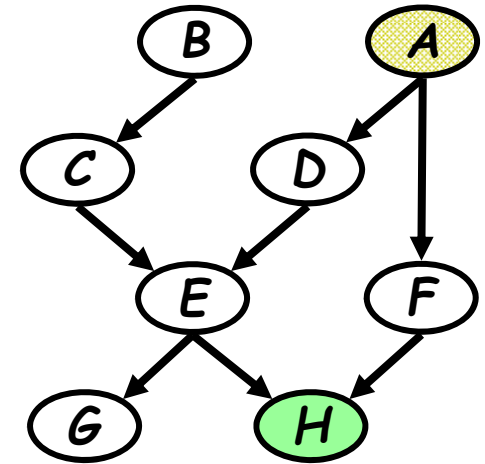
$$\begin{aligned} &P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f) \\ \Rightarrow &P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)m_h(e, f) \\ \Rightarrow &P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)m_h(e, f) \\ \Rightarrow &P(a)P(b)P(c | b)P(d | a)\underline{P(e | c, d)}m_f(a, e) \end{aligned}$$

- Step 4: Eliminate  $E$

- compute

$$m_e(a, c, d) = \sum_e p(e | c, d)m_f(a, e)$$

$$\Rightarrow P(a)P(b)P(c | b)P(d | a)\underline{m_e(a, c, d)}$$





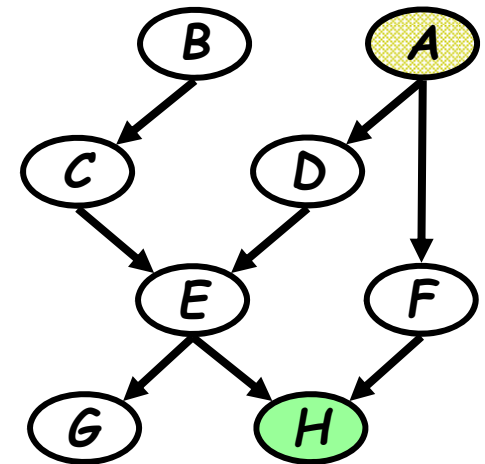


# Example: Variable Elimination

- Query:  $P(B | h)$ 
  - Need to eliminate:  $B, C, D$

- Initial factors:

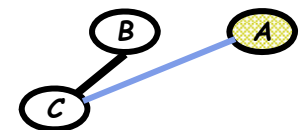
$$\begin{aligned}
 &P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)P(g|e)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)P(e|c,d)m_f(a,e) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)m_e(a,c,d)
 \end{aligned}$$



- Step 5: Eliminate  $D$

- compute  $m_d(a,c) = \sum_d p(d|a)m_e(a,c,d)$

$$\Rightarrow P(a)P(b)P(c|d)m_d(a,c)$$





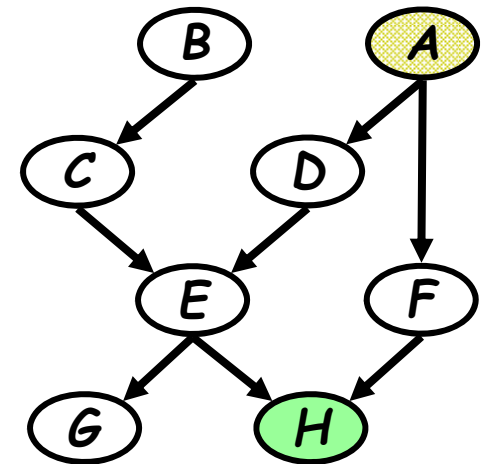
# Example: Variable Elimination

- Query:  $P(B | h)$ 
  - Need to eliminate:  $B, C$
- Initial factors:

$$\begin{aligned} &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)m_h(e, f) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)m_h(e, f) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)m_f(a, e) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)m_e(a, c, d) \\ \Rightarrow &P(a)P(b)P(c | d)m_d(a, c) \end{aligned}$$

- Step 6: **Eliminate  $C$**

- compute 
$$m_c(a, b) = \sum_c p(c | b)m_d(a, c)$$
$$\Rightarrow P(a)P(b)P(c | d)m_d(a, c)$$





# Example: Variable Elimination

- Query:  $P(B | h)$ 
  - Need to eliminate:  $B$
- Initial factors:

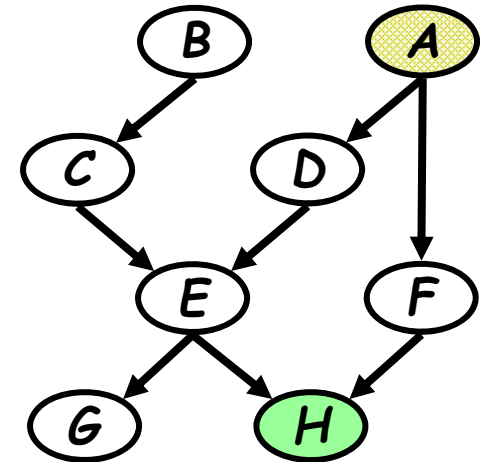
$$\begin{aligned} &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)m_h(e, f) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)m_h(e, f) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)m_f(a, e) \\ \Rightarrow &P(a)P(b)P(c | d)P(d | a)m_e(a, c, d) \\ \Rightarrow &P(a)P(b)P(c | d)m_d(a, c) \\ \Rightarrow &P(a)P(b)m_c(a, b) \end{aligned}$$

- Step 7: **Eliminate  $B$**

- compute

$$m_b(a) = \sum_b p(b)m_c(a, b)$$

$$\Rightarrow \underline{P(a)m_b(a)}$$



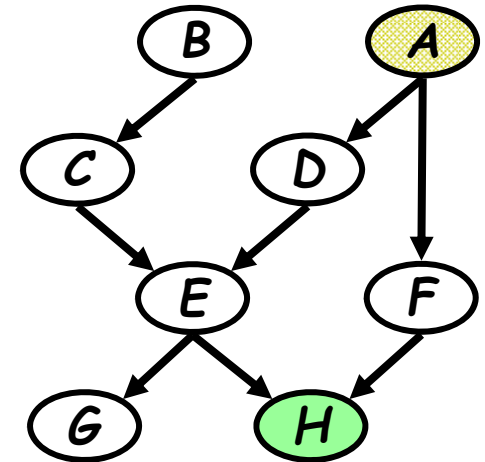


# Example: Variable Elimination

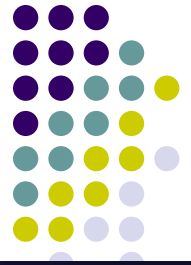
- Query:  $P(B | h)$ 
  - Need to eliminate:  $B$
- Initial factors:

$$\begin{aligned} & P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f) \\ \Rightarrow & P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)m_h(e, f) \\ \Rightarrow & P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)m_h(e, f) \\ \Rightarrow & P(a)P(b)P(c | d)P(d | a)P(e | c, d)m_f(a, e) \\ \Rightarrow & P(a)P(b)P(c | d)P(d | a)m_e(a, c, d) \\ \Rightarrow & P(a)P(b)P(c | d)m_d(a, c) \\ \Rightarrow & P(a)P(b)m_c(a, b) \\ \Rightarrow & P(a)m_b(a) \end{aligned}$$

- Step 8: **Wrap-up**  $p(a, \tilde{h}) = p(a)m_b(a), \quad p(\tilde{h}) = \sum_a p(a)m_b(a)$   
 $\Rightarrow P(a | \tilde{h}) = \frac{p(a)m_b(a)}{\sum_a p(a)m_b(a)}$



# Complexity of variable elimination



- Suppose in one elimination step we compute

$$m_x(y_1, \dots, y_k) = \sum_x m'_x(x, y_1, \dots, y_k)$$
$$m'_x(x, y_1, \dots, y_k) = \prod_{i=1}^k m_i(x, \mathbf{y}_{C_i})$$

This requires

- $k \cdot |\text{Val}(X)| \cdot \prod_i |\text{Val}(\mathbf{y}_{C_i})|$  **multiplications**
  - For each value of  $x, y_1, \dots, y_k$ , we do  $k$  multiplications
- $|\text{Val}(X)| \cdot \prod_i |\text{Val}(\mathbf{y}_{C_i})|$  **additions**
  - For each value of  $y_1, \dots, y_k$ , we do  $|\text{Val}(X)|$  additions

Complexity is **exponential** in number of variables in the intermediate factor

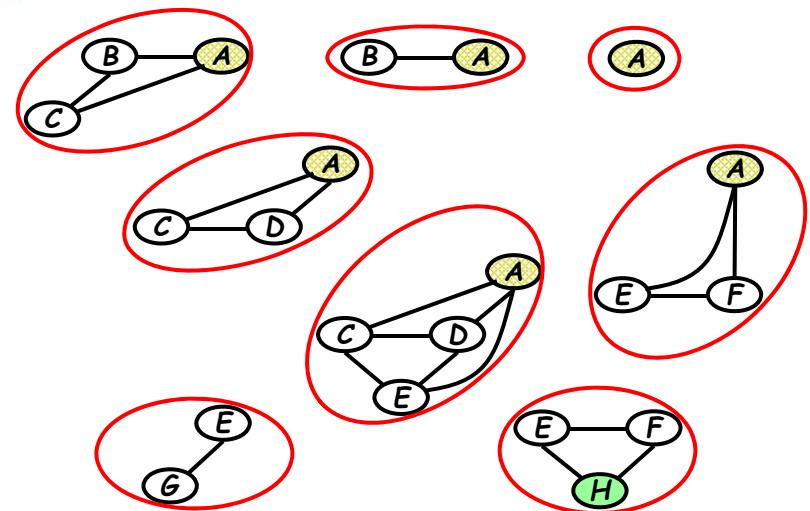


# Elimination Clique

- Induced dependency during marginalization is captured in elimination cliques
  - Summation  $\leftrightarrow$  elimination
  - Intermediate term  $\leftrightarrow$  elimination clique

$$\begin{aligned}
 & P(a)P(b)P(c|b)P(d|a)P(e|c, d)P(f|a)P(g|e)P(h|e, f) \\
 \Rightarrow & P(a)P(b)P(c|b)P(d|a)P(e|c, d)P(f|a)P(g|e)\phi_h(e, f) \\
 \Rightarrow & P(a)P(b)P(c|b)P(d|a)P(e|c, d)P(f|a)\phi_g(e)\phi_h(e, f) \\
 \Rightarrow & P(a)P(b)P(c|b)P(d|a)P(e|c, d)\phi_f(a, e) \\
 \Rightarrow & P(a)P(b)P(c|b)P(d|a)\phi_e(a, c, d) \\
 \Rightarrow & P(a)P(b)P(c|b)\phi_d(a, c) \\
 \Rightarrow & P(a)P(b)\phi_c(a, b) \\
 \Rightarrow & P(a)\phi_b(a) \\
 \Rightarrow & \phi(a)
 \end{aligned}$$

- Can this lead to an generic inference algorithm?

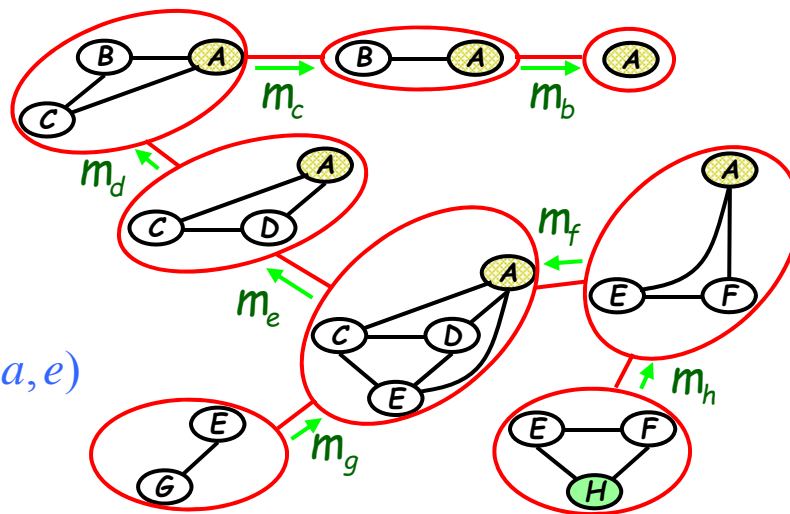


# From Elimination to Message Passing



- Elimination  $\equiv$  message passing on a **clique tree**

$$m_e(a, c, d) = \sum_e p(e | c, d) m_g(e) m_f(a, e)$$

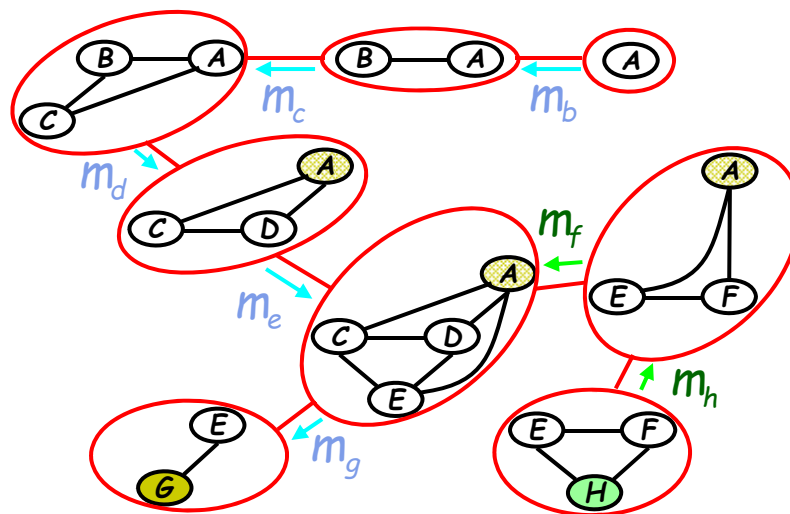


- Messages can be reused

# From Elimination to Message Passing



- Elimination  $\equiv$  message passing on a **clique tree**
  - Another query ...



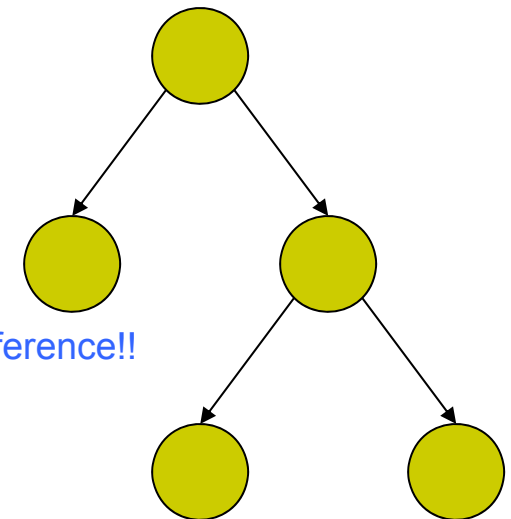
- Messages  $m_f$  and  $m_h$  are reused, others need to be recomputed



# From elimination to message passing

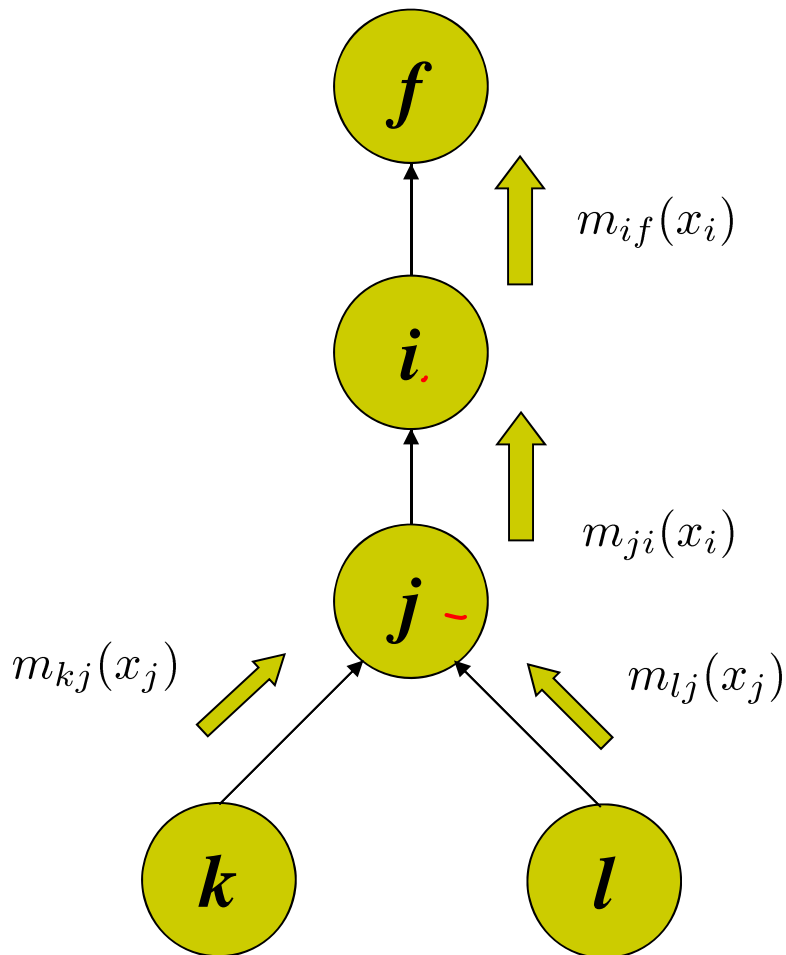


- Recall **ELIMINATION** algorithm:
    - Choose an ordering  $\mathcal{Z}$  in which query node  $f$  is the final node
    - Place all potentials on an active list
    - Eliminate node  $i$  by removing all potentials containing  $i$ , take sum/product over  $x_i$ .
    - Place the resultant factor back on the list
  
  - For a **TREE** graph:
    - Choose query node  $f$  as the root of the tree
    - View tree as a directed tree with edges pointing towards from  $f$
    - Elimination ordering based on depth-first traversal
    - Elimination of each node can be considered as **message-passing (or Belief Propagation)** directly along tree branches, rather than on some transformed graphs
- thus, we can use the tree itself as a data-structure to do general inference!!





# Message passing for trees



Let  $m_{ij}(x_i)$  denote the factor resulting from eliminating variables from below up to  $i$ , which is a function of  $x_i$ :

$$m_{ji}(x_i) = \sum_{x_j} \left( \psi(x_j) \psi(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}(x_j) \right)$$

This is reminiscent of a **message** sent from  $j$  to  $i$ .

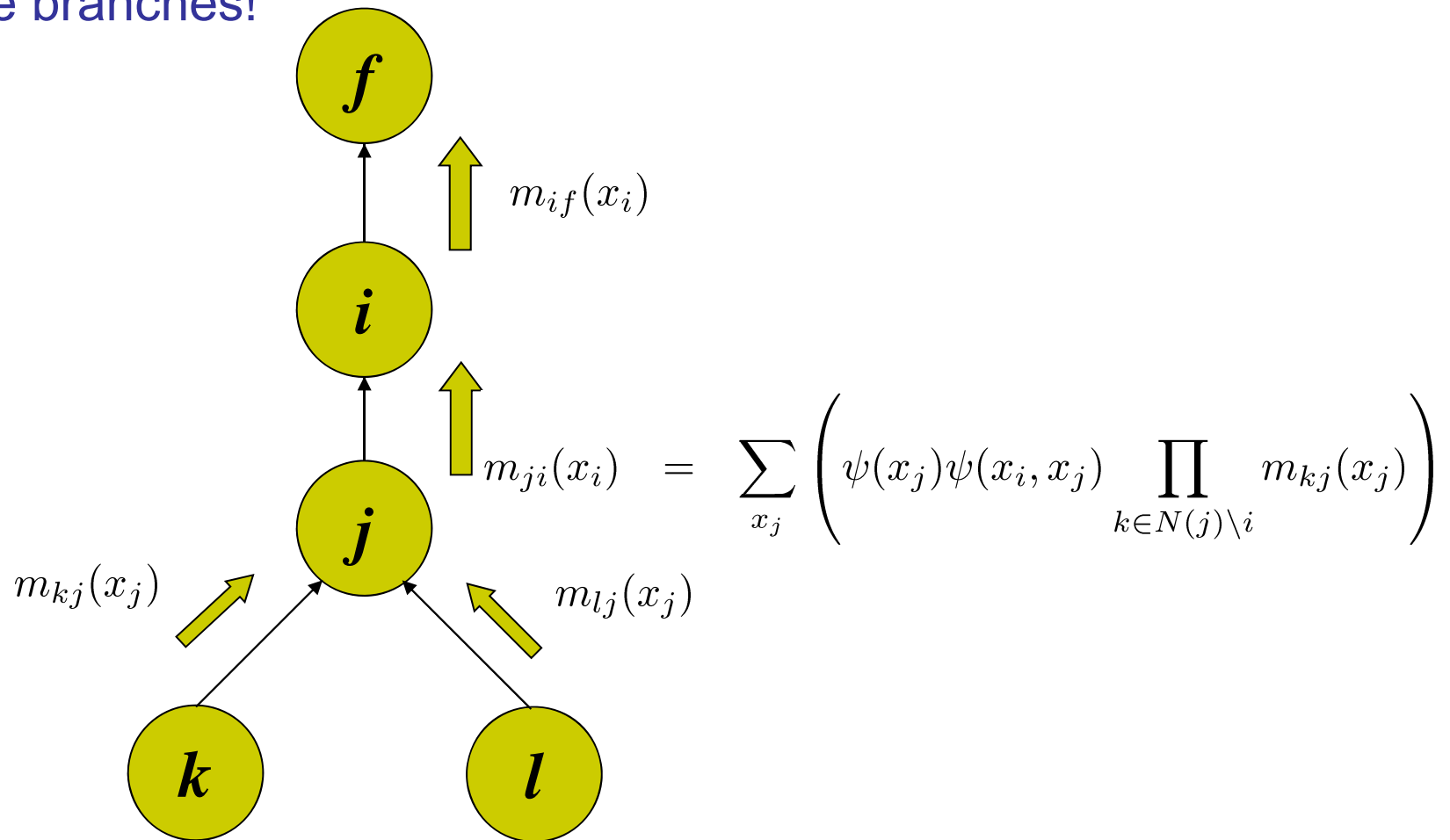
$$m_{ji}(x_i) = \sum_{x_j} \left( \psi(x_j) \psi(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}(x_j) \right)$$

$$p(x_f) \propto \psi(x_f) \prod_{e \in N(f)} m_{ef}(x_f)$$

$m_{ij}(x_i)$  represents a "belief" of  $x_i$  from  $x_j$ !



- Elimination on trees is equivalent to message passing along tree branches!

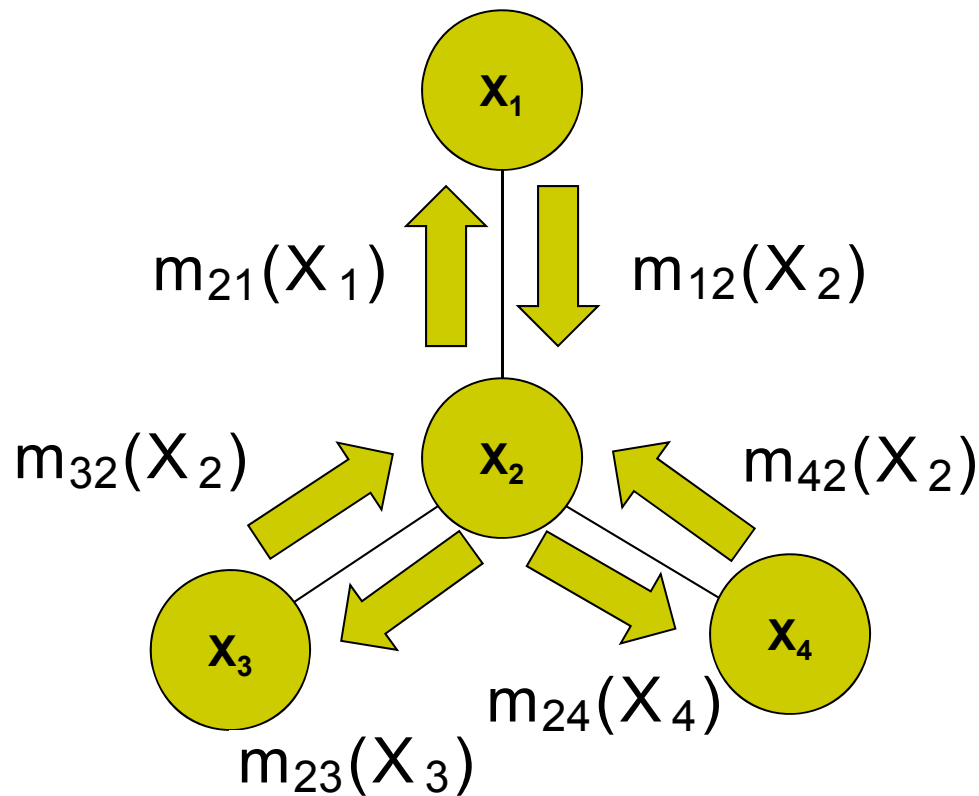


$$= \sum_{x_j} \left( \psi(x_j) \psi(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}(x_j) \right)$$



# The message passing protocol:

- A two-pass algorithm:



# Belief Propagation (SP-algorithm): Sequential implementation



```

SUM-PRODUCT( $\mathcal{T}, E$ )
  EVIDENCE( $E$ )
   $f = \text{CHOOSEROOT}(\mathcal{V})$ 
  for  $e \in \mathcal{N}(f)$ 
    COLLECT( $f, e$ )
  for  $e \in \mathcal{N}(f)$ 
    DISTRIBUTE( $f, e$ )
  for  $i \in \mathcal{V}$ 
    COMPUTEMARGINAL( $i$ )
  
```

```

EVIDENCE( $E$ )
  for  $i \in E$ 
     $\psi^E(x_i) = \psi(x_i)\delta(x_i, \bar{x}_i)$ 
  for  $i \notin E$ 
     $\psi^E(x_i) = \psi(x_i)$ 
  
```

```

COLLECT( $i, j$ )
  for  $k \in \mathcal{N}(j) \setminus i$ 
    COLLECT( $j, k$ )
  SENDMESSAGE( $j, i$ )
  
```

```

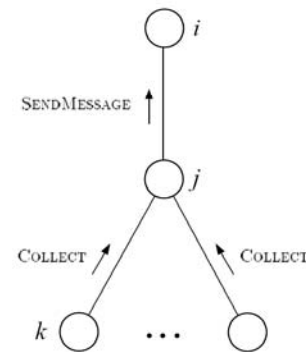
DISTRIBUTE( $i, j$ )
  SENDMESSAGE( $i, j$ )
  for  $k \in \mathcal{N}(j) \setminus i$ 
    DISTRIBUTE( $j, k$ )
  
```

```

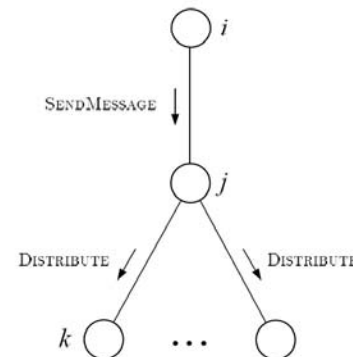
SENDMESSAGE( $j, i$ )
 $m_{ji}(x_i) = \sum_{x_j} (\psi^E(x_j)\psi(x_i, x_j)) \prod_{k \in \mathcal{N}(j) \setminus i} m_{kj}(x_j)$ 
  
```

```

COMPUTEMARGINAL( $i$ )
 $p(x_i) \propto \psi^E(x_i) \prod_{j \in \mathcal{N}(i)} m_{ji}(x_i)$ 
  
```



1





# Inference on general GM

---

- Now, what if the GM is not a tree-like graph?
- Can we still directly run message message-passing protocol along its edges?
- For non-trees, we do not have the guarantee that message-passing will be consistent!
- Then what?
  - Construct a graph data-structure from  $P$  that has a tree structure, and run message-passing on it!

→ Junction tree algorithm



# Summary: Exact Inference

---

- The simple Eliminate algorithm captures the key algorithmic Operation underlying probabilistic inference:  
--- That of taking a sum over product of potential functions
- The computational complexity of the Eliminate algorithm can be reduced to purely graph-theoretic considerations.
- This graph interpretation will also provide hints about how to design improved inference algorithms
- What can we say about the overall computational complexity of the algorithm? In particular, how can we control the "size" of the summands that appear in the sequence of summation operation.

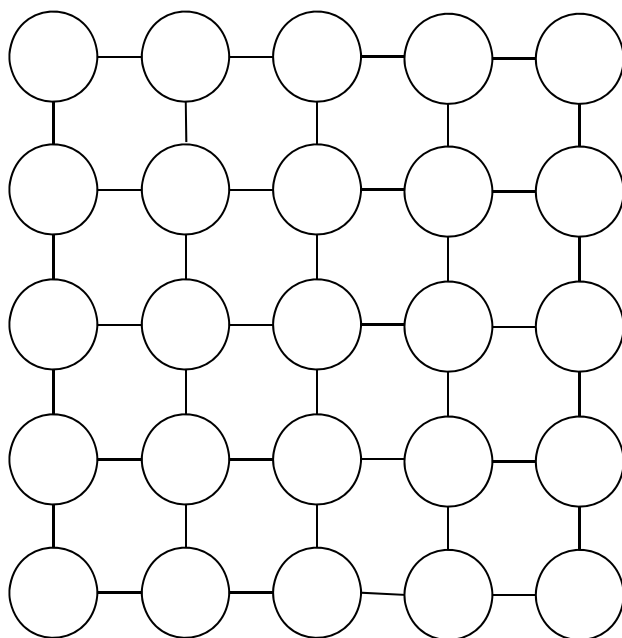


# Approaches to inference

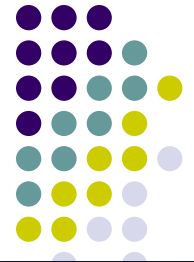
---

- Exact inference algorithms
  - The elimination algorithm
  - Belief propagation
  - The junction tree algorithms (but will not cover in detail here)
- Approximate inference techniques
  - Variational algorithms
  - Stochastic simulation / sampling methods
  - Markov chain Monte Carlo methods





Ising/Potts model



# Monte Carlo methods

- Draw random samples from the desired distribution
- Yield a stochastic representation of a complex distribution
  - marginals and other expectations can be approximated using **sample-based averages**

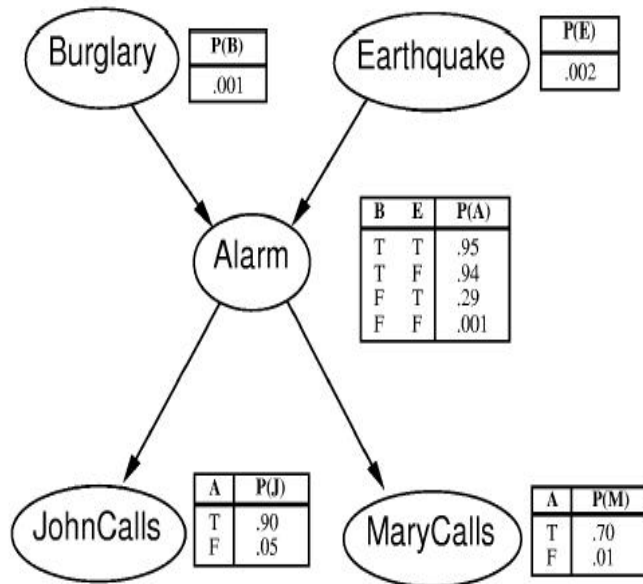
$$E[f(x)] = \frac{1}{N} \sum_{t=1}^N f(x^{(t)})$$

- **Asymptotically** exact and easy to apply to arbitrary models
- Challenges:
  - how to draw samples from a given dist. (not all distributions can be trivially sampled)?
  - how to make better use of the samples (not all sample are useful, or eqally useful, see an example later)?
  - how to know we've sampled enough?



# Example: naive sampling

- Construct samples according to probabilities given in a BN.



E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E1	B0	A1	M1	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0

**Alarm example: (Choose the right sampling sequence)**

1) Sampling:  $P(B) = \langle 0.001, 0.999 \rangle$  suppose it is false, B0. Same for E0.  $P(A|B0, E0) = \langle 0.001, 0.999 \rangle$  suppose it is false...

2) Frequency counting: In the samples right,  $P(J|A0) = P(J, A0) / P(A0) = \langle 1/9, 8/9 \rangle$ .



# Example: naive sampling

- Construct samples according to probabilities given in a BN.

Alarm example: (Choose the right sampling sequence)

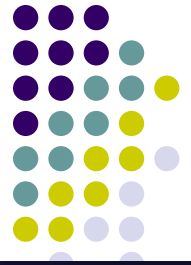
3) what if we want to compute  $P(J|A1)$  ?  
we have only one sample ...  
 $P(J|A1)=P(J,A1)/P(A1)=\langle 0, 1 \rangle$ .

4) what if we want to compute  $P(J|B1)$  ?  
No such sample available!  
 $P(J|A1)=P(J,B1)/P(B1)$  can not be defined.

For a model with hundreds or more variables, rare events will be very hard to garner enough samples even after a long time or sampling ...

E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E1	B0	A1	M1	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0

# Markov chain Monte Carlo (MCMC)



- Construct a Markov chain whose stationary distribution is the target density  $= P(X|e)$ .
- Run for  $T$  samples (burn-in time) until the chain converges/mixes/reaches stationary distribution.
- Then collect  $M$  (correlated) samples  $x_m$  .
- Key issues:
  - Designing proposals so that the chain mixes rapidly.
  - Diagnosing convergence.

# Markov Chains



- **Definition:**

- Given an n-dimensional state space
- Random vector  $\mathbf{X} = (x_1, \dots, x_n)$
- $\mathbf{x}^{(t)} = \mathbf{x}$  at time-step t
- $\mathbf{x}^{(t)}$  transitions to  $\mathbf{x}^{(t+1)}$  with prob

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) = T(\mathbf{x}^{(t)} \rightarrow \mathbf{x}^{(t+1)})$$

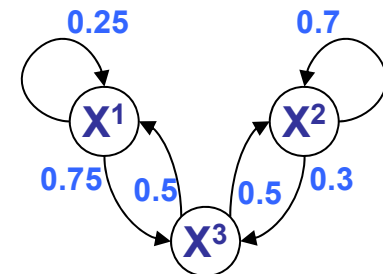
- **Homogenous:** chain determined by state  $\mathbf{x}^{(0)}$ , fixed *transition kernel* Q (rows sum to 1)

- **Equilibrium:**  $\pi(\mathbf{x})$  is a *stationary (equilibrium) distribution* if

$$\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}) Q(\mathbf{x} \rightarrow \mathbf{x}').$$

i.e., is a left eigenvector of the transition matrix  $\pi^T = \pi^T Q$ .

$$(0.2 \quad 0.5 \quad 0.3) = (0.2 \quad 0.5 \quad 0.3) \begin{pmatrix} 0.25 & 0 & 0.75 \\ 0 & 0.7 & 0.3 \\ 0.5 & 0.5 & 0 \end{pmatrix}$$



# Gibbs sampling

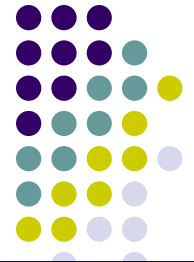


- The transition matrix updates each node one at a time using the following proposal:

$$Q((\mathbf{x}_i, \mathbf{x}_{-i}) \rightarrow (\mathbf{x}_i', \mathbf{x}_{-i})) = p(\mathbf{x}_i' | \mathbf{x}_{-i})$$

- It is efficient since  $p(\mathbf{x}_i' | \mathbf{x}_{-i})$  only depends on the values in  $X_i$ 's Markov blanket

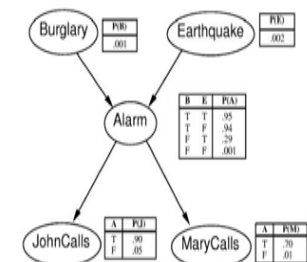
# Gibbs sampling



- Gibbs sampling is an MCMC algorithm that is especially appropriate for inference in graphical models.

- The procedure

- we have variable set  $\mathbf{X} = \{x_1, x_2, x_3, \dots, x_N\}$  for a GM
- at each step one of the variables  $X_i$  is selected (at random or according to some fixed sequences), denote the remaining variables as  $\mathbf{X}_{-i}$ , and its current value as  $\mathbf{x}_{-i}^{(t-1)}$ 
  - Using the "alarm network" as an example, say at time  $t$  we choose  $X_E$ , and we denote the current value assignments of the remaining variables,  $\mathbf{X}_{-E}$ , obtained from previous samples, as  $\mathbf{x}_{-E}^{(t-1)} = \{x_B^{(t-1)}, x_A^{(t-1)}, x_J^{(t-1)}, x_M^{(t-1)}\}$
- the conditional distribution  $p(X_i | \mathbf{x}_{-i}^{(t-1)})$  is computed
- a value  $x_i^{(t)}$  is sampled from this distribution
- the sample  $x_i^{(t)}$  replaces the previous sampled value of  $X_i$  in  $\mathbf{X}$ .
  - i.e.,  $\mathbf{x}^{(t)} = \mathbf{x}_{-E}^{(t-1)} \cup x_E^{(t)}$







# Markov Blanket

## ● Markov Blanket in BN

- A variable is independent from others, given its parents, children and children's parents (d-separation).

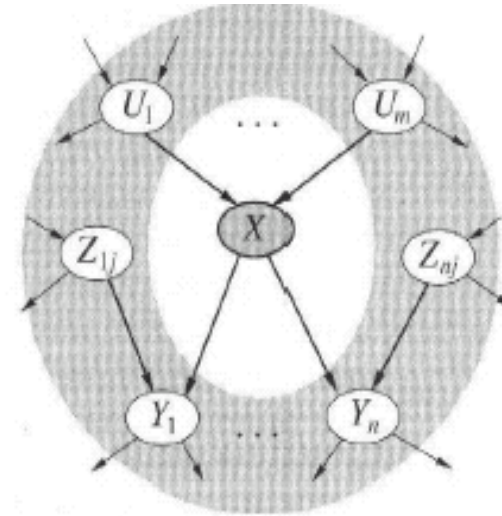
## ● MB in MRF

- A variable is independent all its non-neighbors, given all its direct neighbors.

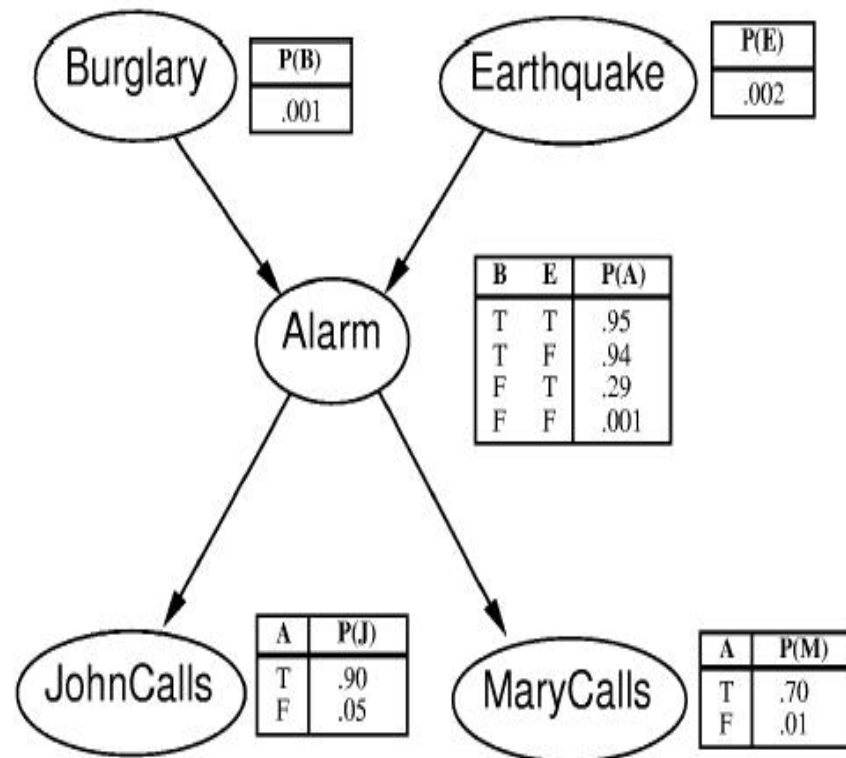
$$\Rightarrow p(X_i | X_{\setminus i}) = p(X_i | \text{MB}(X_i))$$

## ● Gibbs sampling

- Every step, choose one variable and sample it by  $P(X|\text{MB}(X))$  based on previous sample.



# Gibbs sampling of the alarm network

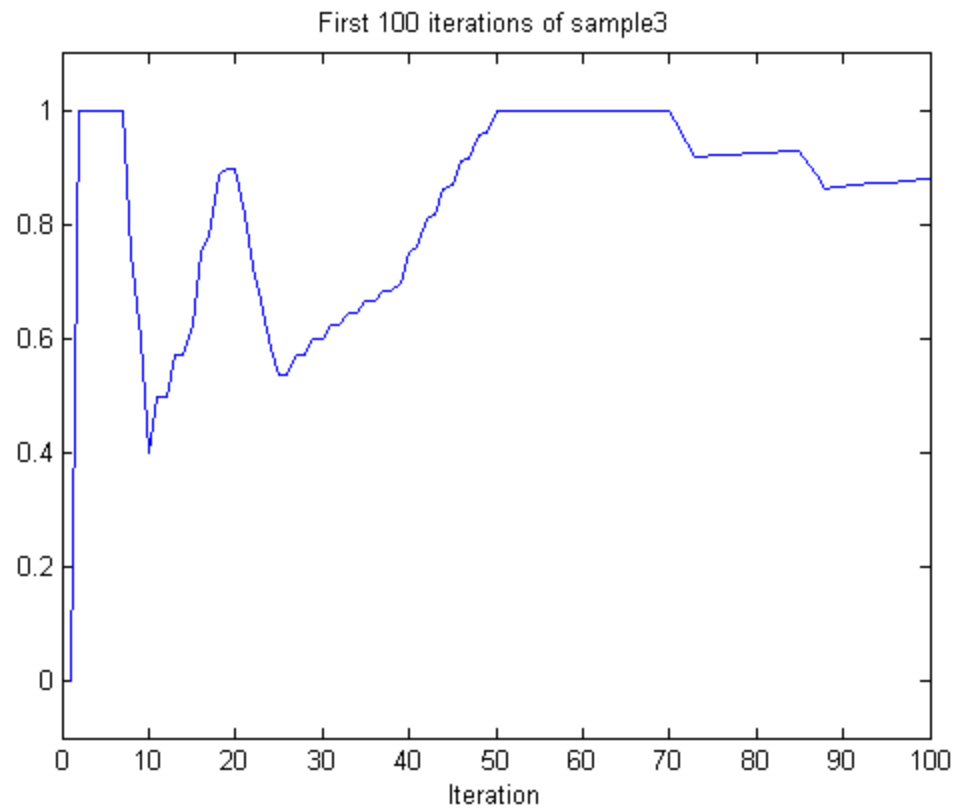
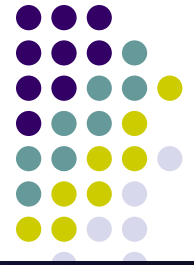


**$MB(A) = \{B, E, J, M\}$**

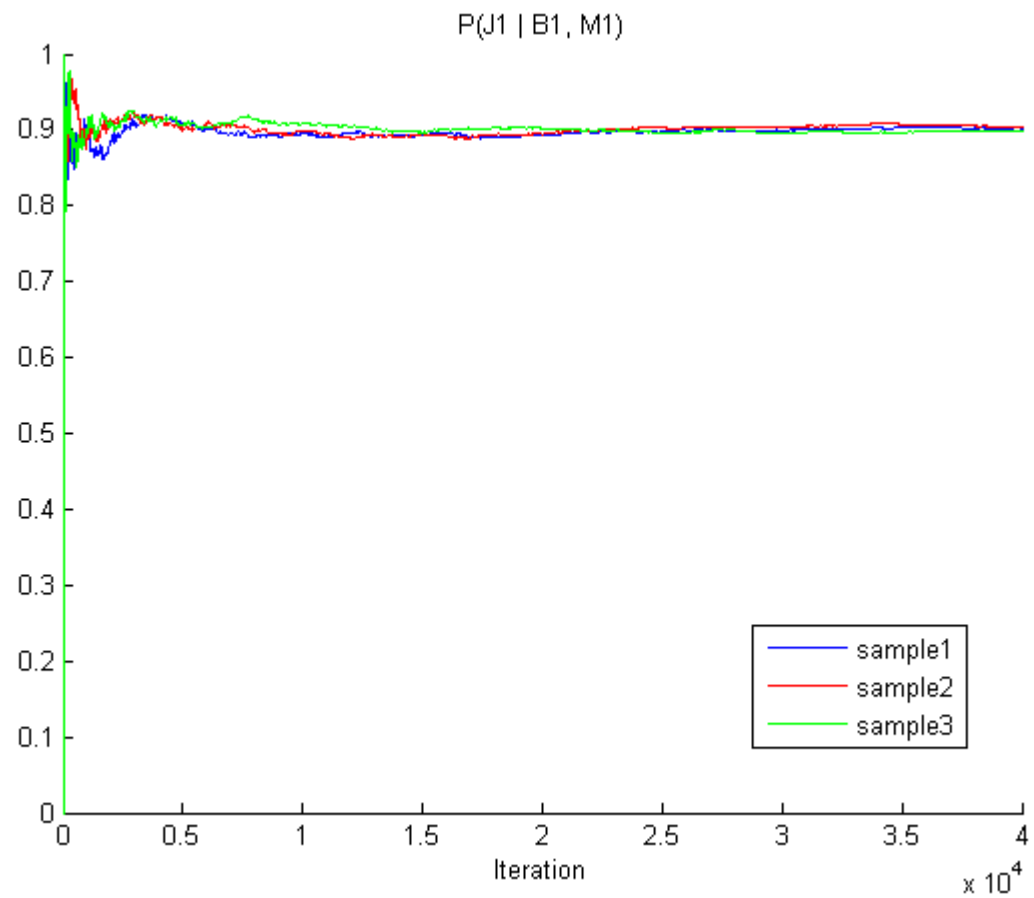
**$MB(E) = \{A, B\}$**

- To calculate  $P(J|B1, M1)$
- Choose  $(B1, E0, A1, M1, J1)$  as a start
- Evidences are  $B1, M1$ , variables are  $A, E, J$ .
- Choose next variable as  $A$
- Sample  $A$  by  $P(A|MB(A)) = P(A|B1, E0, M1, J1)$  suppose to be false.
- $(B1, E0, A0, M1, J1)$
- Choose next random variable as  $E$ , sample  $E \sim P(E|B1, A0)$
- ...

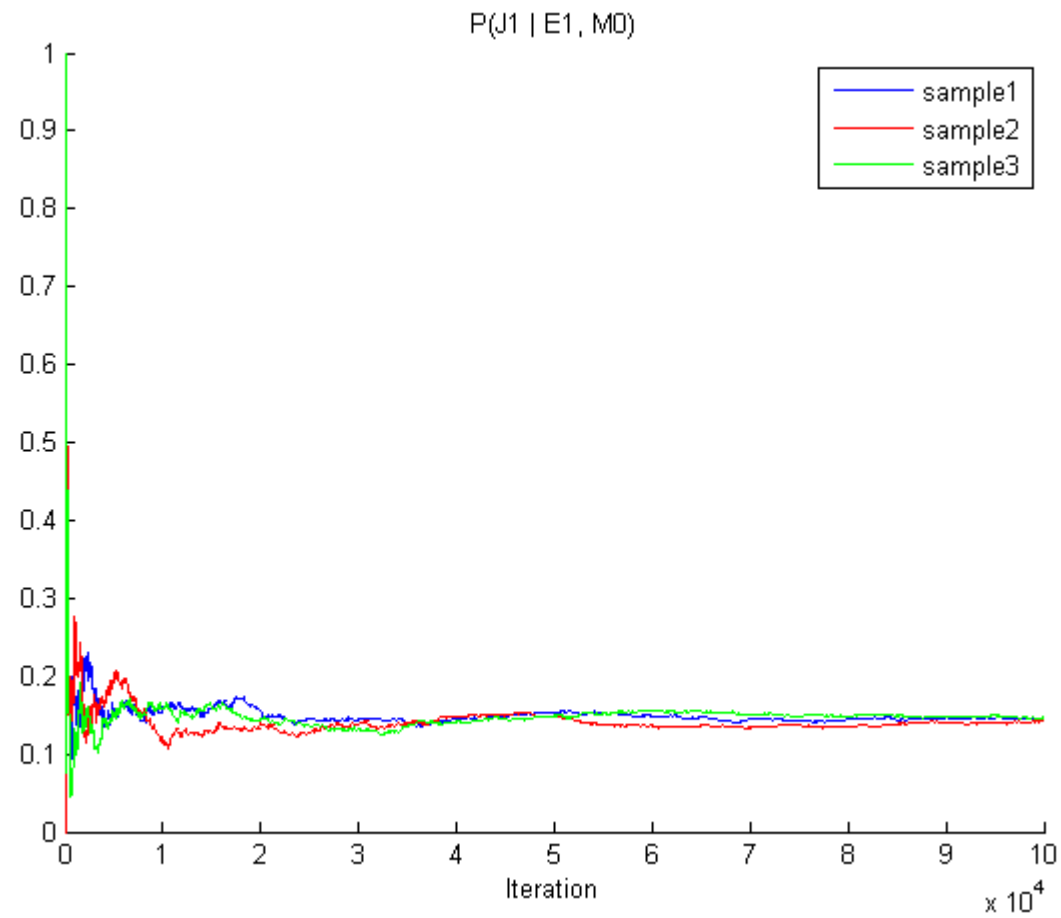
# Example



# Example:



# Example



# Example



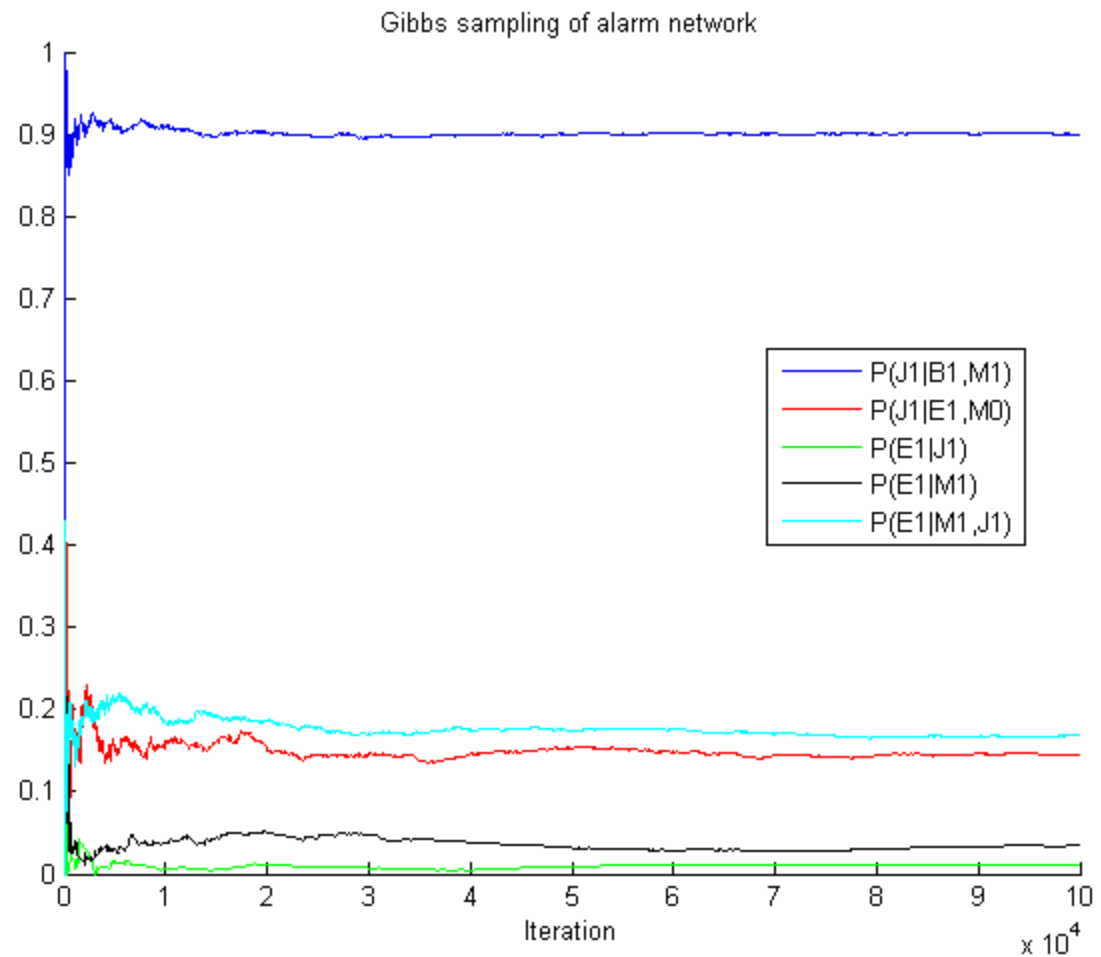
$$P(J1 | B1, M1) = 0.90$$

$$P(J1 | E1, M0) = 0.14$$

$$P(E1 | J1) = 0.01$$

$$P(E1 | M1) = 0.04$$

$$P(E1 | M1, J1) = 0.17$$



# The **Art** of simulation



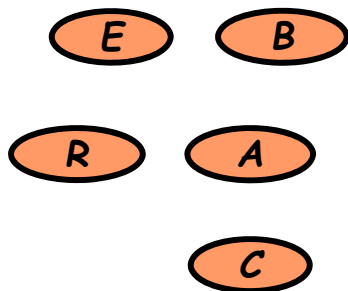
- Run several chains
- Start at over-dispersed points
- Monitor the log lik.
- Monitor the serial correlations
- Monitor acceptance ratios
- Re-parameterize (to get approx. indep.)
- Re-block (Gibbs)
- Collapse (int. over other pars.)
- Run with troubled pars. fixed at reasonable vals.



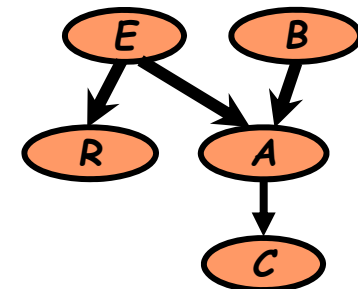
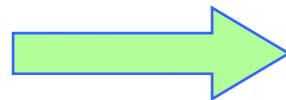
# Learning Graphical Models

## The goal:

Given set of independent samples (*assignments of random variables*), find the *best* (the most likely?) Bayesian Network (both DAG and CPDs)



(B,E,A,C,R)=(T,F,F,T,F)  
(B,E,A,C,R)=(T,F,T,T,F)  
.....  
(B,E,A,C,R)=(F,T,T,T,F)



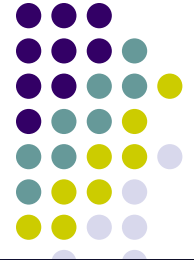
E	B	P(A   E, B)	
e	b	0.9	0.1
e	$\bar{b}$	0.2	0.8
$\bar{e}$	b	0.9	0.1
$\bar{e}$	$\bar{b}$	0.01	0.99



# Learning Graphical Models (cont.)



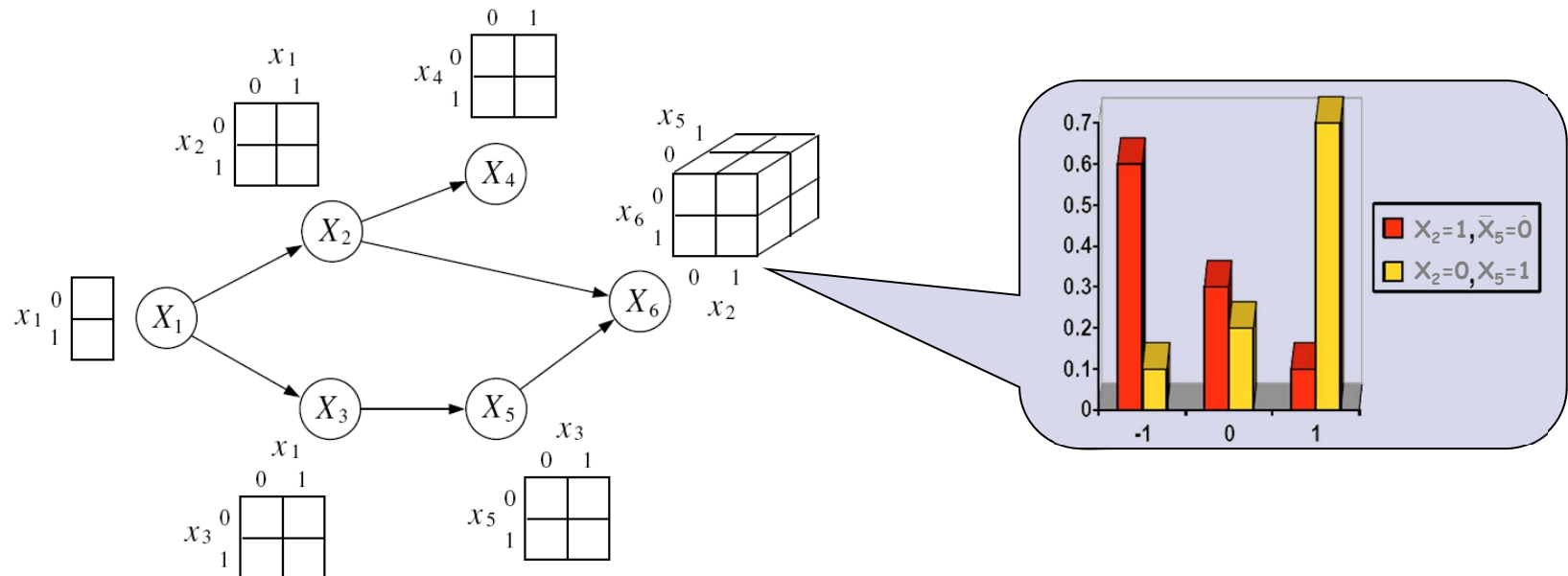
- Scenarios:
  - completely observed GMs
    - directed
    - undirected
  - partially observed GMs
    - directed
    - undirected (an open research topic)
- Estimation principles:
  - Maximal likelihood estimation (MLE)
  - Bayesian estimation
  - Maximal conditional likelihood
  - Maximal "Margin"
- We use **learning** as a name for the process of **estimating the parameters**, and in some cases, the topology of the network, from data.



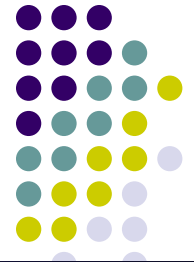
# MLE for general BN parameters

- If we assume the parameters for each CPD are globally independent, and all nodes are **fully observed**, then the log-likelihood function decomposes into a sum of local terms, one per node:

$$\ell(\theta; D) = \log p(D | \theta) = \log \prod_i \left( \prod_n p(x_{n,i} | \mathbf{x}_{n,\pi_i}, \theta_i) \right) = \sum_i \left( \sum_n \log p(x_{n,i} | \mathbf{x}_{n,\pi_i}, \theta_i) \right)$$



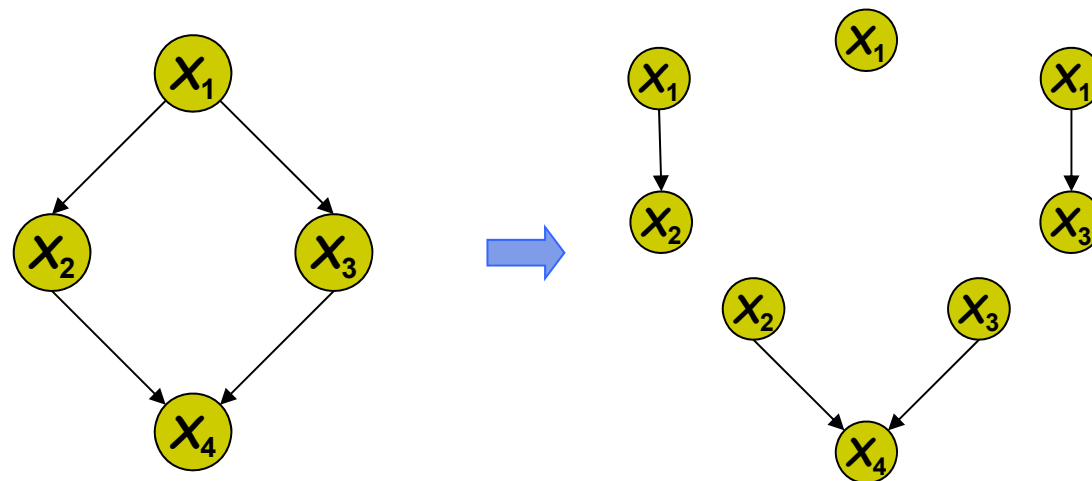
# Example: decomposable likelihood of a directed model



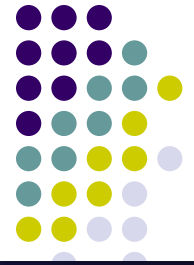
- Consider the distribution defined by the directed acyclic GM:

$$p(x|\theta) = p(x_1|\theta_1)p(x_2|x_1,\theta_2)p(x_3|x_1,\theta_3)p(x_4|x_2,x_3,\theta_4)$$

- This is exactly like learning four separate small BNs, each of which consists of a node and its parents.



# E.g.: MLE for BNs with tabular CPDs



- Assume each CPD is represented as a table (multinomial) where

$$\theta_{ijk} \stackrel{\text{def}}{=} p(X_i = j \mid X_{\pi_i} = k)$$

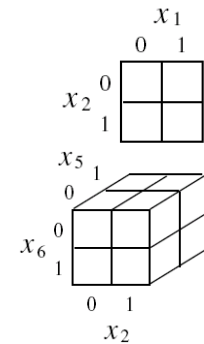
- Note that in case of multiple parents,  $\mathbf{X}_{\pi_i}$  will have a composite state, and the CPD will be a high-dimensional table
- The sufficient statistics are counts of family configurations

$$n_{ijk} \stackrel{\text{def}}{=} \sum_n x_{n,i}^j x_{n,\pi_i}^k$$

- The log-likelihood is  $\ell(\theta; \mathcal{D}) = \log \prod_{i,j,k} \theta_{ijk}^{n_{ijk}} = \sum_{i,j,k} n_{ijk} \log \theta_{ijk}$

- Using a Lagrange multiplier to enforce  $\sum_j \theta_{ijk} = 1$ , we get:

$$\theta_{ijk}^{ML} = \frac{n_{ijk}}{\sum_{i,j',k} n_{ij'k}}$$



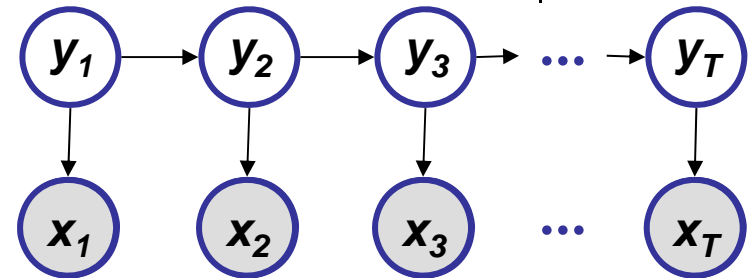


# Definition of HMM

- Transition probabilities between any two states

$$p(y_t^j = \mathbf{1} | y_{t-1}^i = \mathbf{1}) = a_{i,j},$$

or  $p(y_t | y_{t-1}^i = \mathbf{1}) \sim \text{Multinomial}(a_{i,1}, a_{i,2}, \dots, a_{i,M}), \forall i \in \mathcal{I}.$



- Start probabilities

$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, \dots, \pi_M).$$

- Emission probabilities associated with each state

$$p(x_t | y_t^i = \mathbf{1}) \sim \text{Multinomial}(b_{i,1}, b_{i,2}, \dots, b_{i,K}), \forall i \in \mathcal{I}.$$

or in general:  $p(x_t | y_t^i = \mathbf{1}) \sim f(\cdot | \theta_i), \forall i \in \mathcal{I}.$



# Supervised ML estimation

- Given  $\mathbf{x} = x_1 \dots x_N$  for which the true state path  $\mathbf{y} = y_1 \dots y_N$  is known,

$$\ell(\theta; \mathbf{x}, \mathbf{y}) = \log p(\mathbf{x}, \mathbf{y}) = \log \prod_n \left( p(y_{n,1}) \prod_{t=2}^T p(y_{n,t} | y_{n,t-1}) \prod_{t=1}^T p(x_{n,t} | x_{n,t}) \right)$$

- Define:

$A_{ij}$  = # times state transition  $i \rightarrow j$  occurs in  $\mathbf{y}$

$B_{ik}$  = # times state  $i$  in  $\mathbf{y}$  emits  $k$  in  $\mathbf{x}$

- We can show that the **maximum likelihood** parameters  $\theta$  are:

$$a_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^T y_{n,t-1}^i y_{n,t}^j}{\sum_n \sum_{t=2}^T y_{n,t-1}^i} = \frac{A_{ij}}{\sum_{j'} A_{ij'}}$$

$$b_{ik}^{ML} = \frac{\#(i \rightarrow k)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=1}^T y_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^T y_{n,t}^i} = \frac{B_{ik}}{\sum_{k'} B_{ik'}}$$

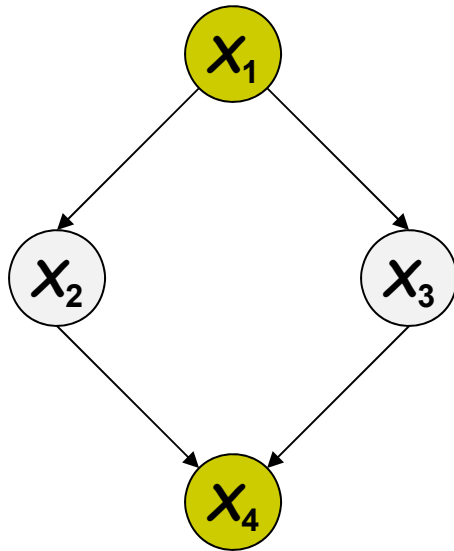
- If  $\mathbf{y}$  is continuous, we can treat  $\{(x_{n,t}, y_{n,t}) : t = 1:T, n = 1:N\}$  as  $N \times T$  observations of, e.g., a Gaussian, and apply learning rules for Gaussian ...

# What if some nodes are not observed?



- Consider the distribution defined by the directed acyclic GM:

$$p(x | \theta) = p(x_1 | \theta_1) p(x_2 | x_1, \theta_1) p(x_3 | x_1, \theta_3) p(x_4 | x_2, x_3, \theta_1)$$



- Need to compute  $p(x_H | x_V) \rightarrow$  inference

# MLE for BNs with tabular CPDs



- Assume each CPD is represented as a table (multinomial) where

$$\theta_{ijk} \stackrel{\text{def}}{=} p(X_i = j \mid X_{\pi_i} = k)$$

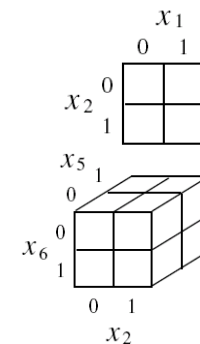
- Note that in case of multiple parents,  $\mathbf{X}_{\pi_i}$  will have a composite state, and the CPD will be a high-dimensional table
- The sufficient statistics are counts of family configurations

$$n_{ijk} \stackrel{\text{def}}{=} \sum_n \langle x_{n,i}^j x_{n,\pi_i}^k \rangle$$

- The log-likelihood is  $\ell(\theta; \mathcal{D}) = \log \prod_{i,j,k} \theta_{ijk}^{n_{ijk}} = \sum_{i,j,k} n_{ijk} \log \theta_{ijk}$

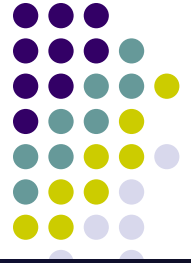
- Using a Lagrange multiplier to enforce  $\sum_j \theta_{ijk} = 1$ , we get:

$$\theta_{ijk}^{ML} = \frac{n_{ijk}}{\sum_{i,j',k} n_{ij'k}}$$





# Summary



- A GMM describes a unique probability distribution  $P$
- Typical tasks:
  - Task 1: How do we answer **queries** about  $P$ ?
    - We use **inference** as a name for the process of computing answers to such queries
  - Task 2: How do we estimate a **plausible model**  $M$  from data  $D$ ?
    - i. We use **learning** as a name for the process of obtaining point estimate of  $M$ .
    - ii. But for *Bayesian*, they seek  $p(M | D)$ , which is actually an **inference** problem.
    - iii. When not all variables are observable, even computing point estimate of  $M$  need to do **inference** to impute the *missing data*.