

CMSC 726

Lecture 20: Inference in Graphical Models

Lise Getoor
November 11, 2010

ACKNOWLEDGEMENTS: The material in this course is a synthesis of materials from many sources, including: Hal Daume III, Mark Drezde, Carlos Guestrin, Andrew Ng, Ben Taskar, **Eric Xing**, and others. I am very grateful for their generous sharing of insights and materials.

Probabilistic Inference

- ▶ We now have compact representations of probability distributions: **Graphical Models**
- ▶ A GM \mathcal{M} describes a unique probability distribution \mathcal{P}
- ▶ How do we answer **queries** about \mathcal{P} ?
- ▶ We use **inference** for the name of the process of computing answers to queries

Query 1: Likelihood

- ▶ Most of the queries one may ask involve **evidence**
 - Evidence e is an assignment of values to a set E variables in the domain
 - Without loss of generality $E = \{ X_{k+1}, \dots, X_n \}$
- ▶ Simplest query: compute probability of evidence

$$P(e) = \sum_{x_1} \dots \sum_{x_k} P(x_1, \dots, x_k, e)$$

- this is often referred to as computing the **likelihood** of e

Query 2: Conditional Probability

- ▶ Often we are interested in the **conditional probability distribution** of a variable given the evidence

$$P(X | e) = \frac{P(X, e)}{P(e)} = \frac{P(X, e)}{\sum_x P(X = x, e)}$$

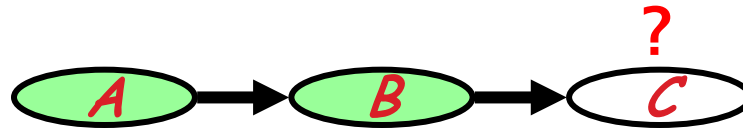
- this is the **a posteriori belief** in X , given evidence e
- ▶ We usually query a subset Y of all domain variables $X = \{Y, Z\}$ and "don't care" about the remaining, Z :

$$P(Y | e) = \sum_z P(Y, Z = z | e)$$

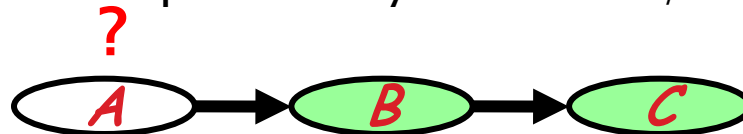
- the process of summing out the "don't care" variables z is called **marginalization**, and the resulting $P(y|e)$ is called a **marginal** probability

Applications of *a posteriori* Belief

- ▶ **Prediction:** what is the probability of an outcome given the starting condition



- the query node is a descendent of the evidence
- ▶ **Diagnosis:** what is the probability of disease/fault given symptoms



- the query node an ancestor of the evidence
- ▶ **Learning** under partial observation
 - fill in the unobserved values under an "EM" setting (more later)
- ▶ The directionality of information flow between variables is *not* restricted by the directionality of the edges in a GM
 - probabilistic inference can combine evidence from all parts of the network

Query 3: Most Probable Assignment

- ▶ In this query we want to find the **most probable joint assignment** (MPA) for *some* variables of interest
- ▶ Such reasoning is usually performed under some given evidence e , and ignoring (the values of) other variables z :

$$\text{MPA}(Y \mid e) = \arg \max_y P(y \mid e) = \arg \max_y \sum_z P(y, z \mid e)$$

- this is the **maximum *a posteriori*** configuration of y .

Applications of MPA

- ▶ Classification
 - find most likely label, given the evidence
- ▶ Explanation
 - what is the most likely scenario, given the evidence

Cautionary note:

- ▶ The MPA of a variable depends on its "context"---the set of variables been jointly queried
- ▶ Example:
 - MPA of X ?
 - MPA of (X, Y) ?

x	y	$P(x,y)$
0	0	0.35
0	1	0.05
1	0	0.3
1	1	0.3

Complexity of Inference

Thm:

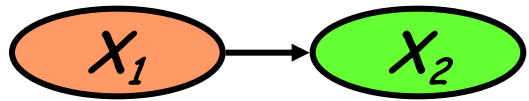
Computing $P(X = x \mid e)$ in a GM is NP-hard

- Hardness does not mean we cannot solve inference
 - It implies that we cannot find a general procedure that works efficiently for arbitrary GMs
 - For particular families of GMs, we can have provably efficient procedures

Approaches to inference

- ▶ Exact inference algorithms
 - The elimination algorithm
 - The junction tree algorithms (not covered in detail)
- ▶ Approximate inference techniques
 - Stochastic simulation / sampling methods
 - Markov chain Monte Carlo methods
 - Variational algorithms (no covered)

Inference in Simple Chains



How do we compute $P(X_2)$?

$$P(x_2) = \sum_{x_1} P(x_1, x_2) = \sum_{x_1} P(x_1) P(x_2 \mid x_1)$$

Inference in Simple Chains (cont.)



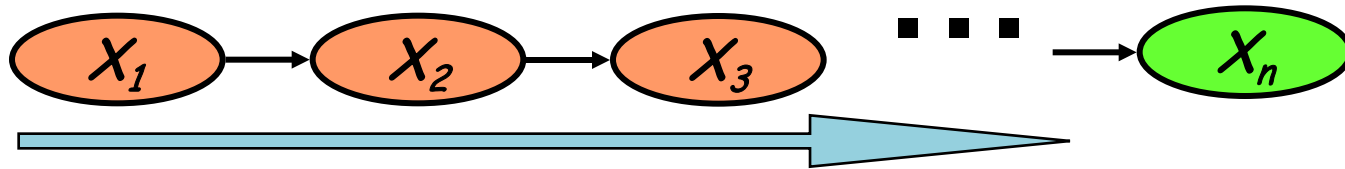
How do we compute $P(X_3)$?

$$P(x_3) = \sum_{x_2} P(x_2, x_3) = \sum_{x_2} P(x_2)P(x_3 | x_2)$$

► we already know how to compute $P(X_2)$...

$$P(x_2) = \sum_{x_1} P(x_1, x_2) = \sum_{x_1} P(x_1)P(x_2 | x_1)$$

Inference in Simple Chains (cont.)



How do we compute $P(X_n)$?

- ▶ Compute $P(X_1)$, $P(X_2)$, $P(X_3)$, ...
- ▶ We compute each term by using the previous one

$$P(x_{i+1}) = \sum_{x_i} P(x_i) P(x_{i+1} | x_i)$$

Complexity:

- Each step costs $O(|Val(X_i)| * |Val(X_{i+1})|)$ operations
- Compare to naïve evaluation, that requires summing over joint values of $n-1$ variables

Elimination in Chains

- ▶ We now try to understand the simple chain example using first-order principles



- ▶ Using definition of probability, we have

$$P(e) = \sum_d \sum_c \sum_b \sum_a P(a, b, c, d, e)$$

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

Elimination in Chains



- ▶ By chain decomposition, we get

$$\begin{aligned} P(e) &= \sum_d \sum_c \sum_b \sum_a P(a, b, c, d, e) \\ &= \sum_d \sum_c \sum_b \sum_a P(a) P(b | a) P(c | b) P(d | c) P(e | d) \end{aligned}$$

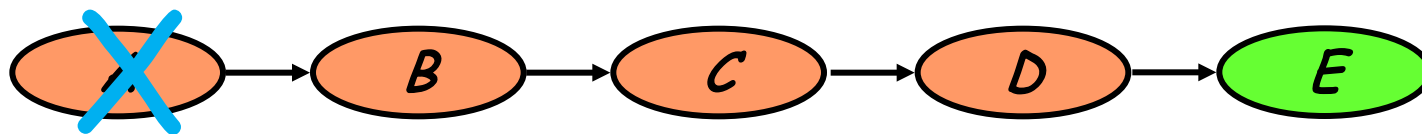
Elimination in Chains



- ▶ Rearranging terms ...

$$\begin{aligned} P(e) &= \sum_d \sum_c \sum_b \sum_a P(a)P(b|a)P(c|b)P(d|c)P(e|d) \\ &= \sum_d \sum_c \sum_b P(c|b)P(d|c)P(e|d) \sum_a P(a)P(b|a) \end{aligned}$$

Elimination in Chains

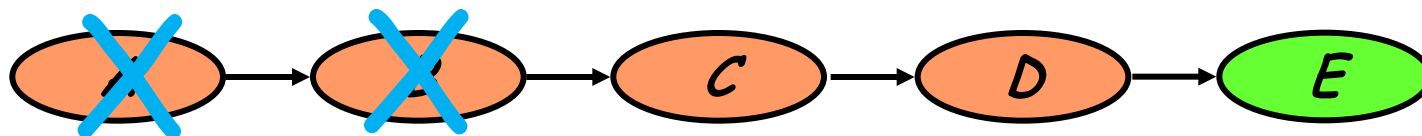


- ▶ Now we can perform innermost summation

$$\begin{aligned} P(e) &= \sum_d \sum_c \sum_b P(c|b)P(d|c)P(e|d) \underbrace{\sum_a P(a)P(b|a)}_{p(b)} \\ &= \sum_d \sum_c \sum_b P(c|b)P(d|c)P(e|d)p(b) \end{aligned}$$

- ▶ This summation, is exactly the first step in the forward iteration we describe before

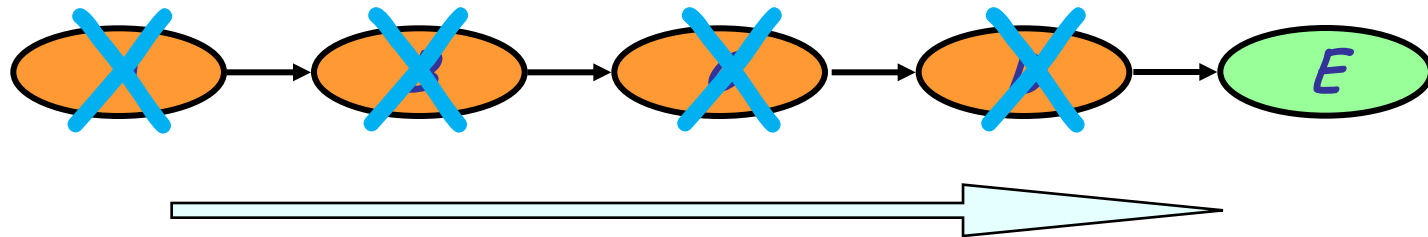
Elimination in Chains



- Rearranging and then summing again, we get

$$\begin{aligned} P(e) &= \sum_d \sum_c \sum_b P(c|b) P(d|c) P(e|d) p(b) \\ &= \sum_d \sum_c P(d|c) P(e|d) \underbrace{\sum_b P(c|b) p(b)}_{p(c)} \\ &= \sum_d \sum_c P(d|c) P(e|d) p(c) \end{aligned}$$

Elimination in Chains



$$P(e) = \sum_d P(e | d) p(d)$$

- ▶ Eliminate nodes one by one all the way to the end, we get

Inference on General GM via Variable Elimination

General idea:

- ▶ Write query in the form

$$P(X_1, \mathbf{e}) = \sum_{x_n} \cdots \sum_{x_3} \sum_{x_2} \prod_i P(x_i \mid pa_i)$$

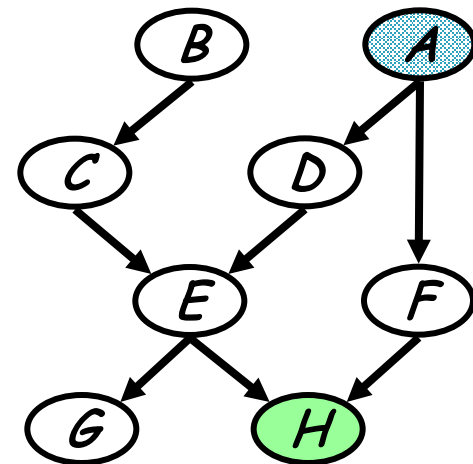
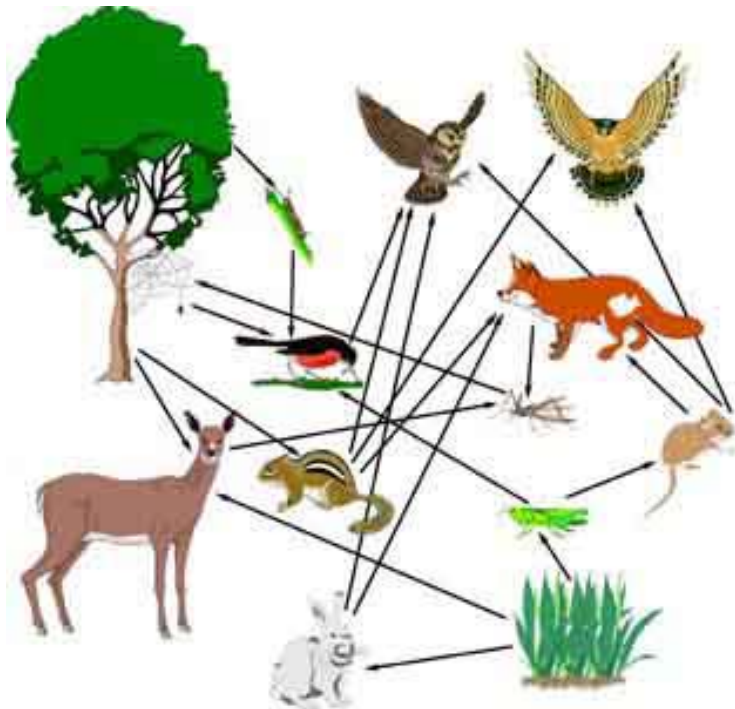
- this suggests an "elimination order" of variables to be marginalized
- ▶ Iteratively
 - Move all irrelevant terms outside of **innermost sum**
 - Perform innermost sum, getting a **new term**
 - Insert the new term into the product

- ▶ wrap-up

$$P(X_1 \mid \mathbf{e}) = \frac{P(X_1, \mathbf{e})}{P(\mathbf{e})}$$

A more complex network

A food web

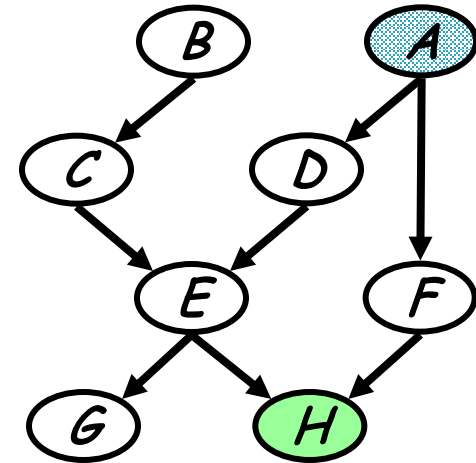


What is the probability that the grass condition is poor given that hawks are leaving?

Example: 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

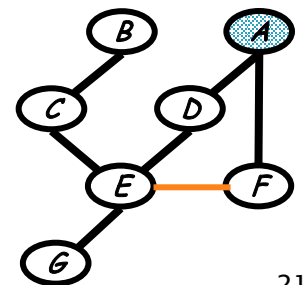


- ▶ Step 1:
 - **Conditioning** (fix the evidence node (i.e., h) on its observed value (i.e., \tilde{h})):

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

- 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:

$$P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f)$$

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

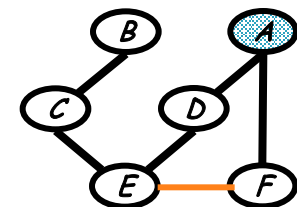
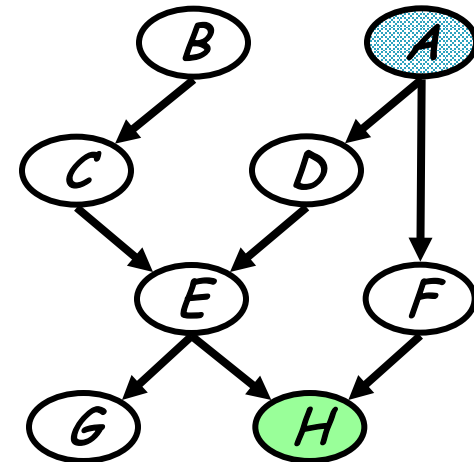
- ▶ Step 2: Eliminate G

- compute

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

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

$$= P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)m_h(e, f)$$



Example: Variable Elimination

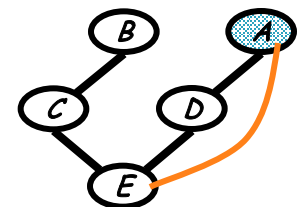
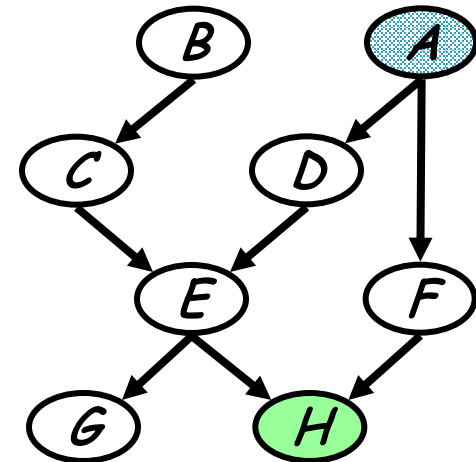
- ▶ Query: $P(B | h)$
 - Need to eliminate: B, C, D, E, F
- ▶ Initial factors:

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

- ▶ Step 3: **Eliminate F**
 - compute

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

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



Example: Variable Elimination

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

- ▶ Initial factors:

$$P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f)$$

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

$$\Rightarrow P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)m_h(e, f)$$

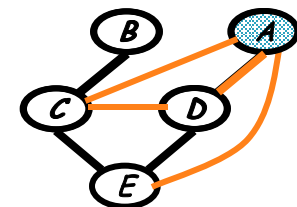
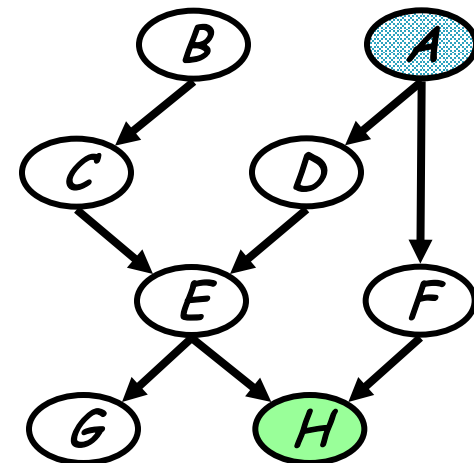
$$\Rightarrow P(a)P(b)P(c | b)P(d | a)\underline{P(e | c, d)m_f(a, e)}$$

- ▶ 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:

$$P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f)$$

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

$$\Rightarrow P(a)P(b)P(c | b)P(d | a)P(e | c, d)P(f | a)m_h(e, f)$$

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

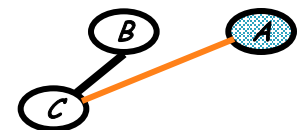
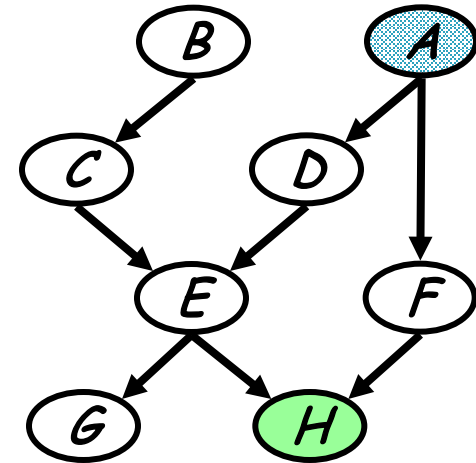
$$\Rightarrow P(a)P(b)P(c | b)P(d | a)m_e(a, c, d)$$

- ▶ 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:

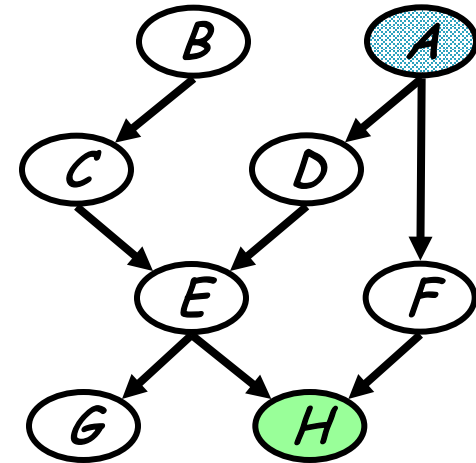
$$\begin{aligned}
 &P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)P(e|c,d)m_f(a,e) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)m_e(a,c,d) \\
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

- ▶ 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:

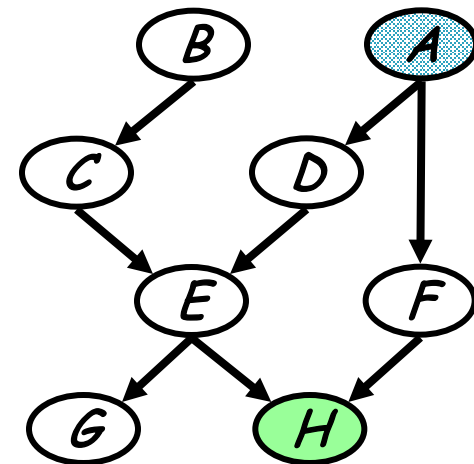
$$\begin{aligned}
 &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)P(h | e, f) \\
 \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)P(g | e)m_h(e, f) \\
 \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)P(f | a)m_h(e, f) \\
 \Rightarrow &P(a)P(b)P(c | d)P(d | a)P(e | c, d)m_f(a, e) \\
 \Rightarrow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

- ▶ 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

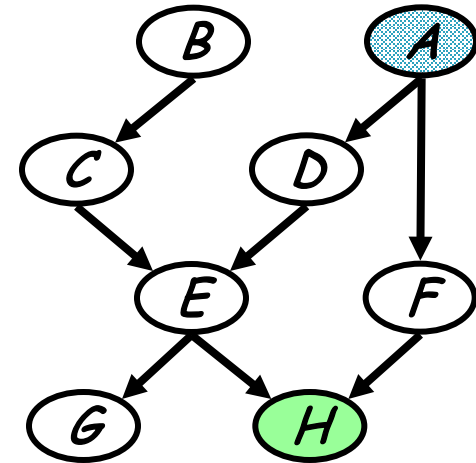
- ▶ Initial factors:

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

- ▶ 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

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

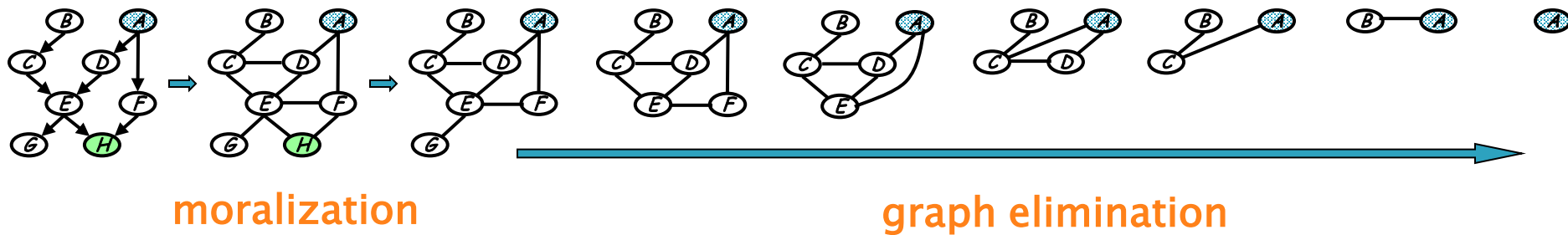
This requires

- ▶ $k \cdot |\text{Val}(X)| \cdot \prod_i |\text{Val}(\mathbf{y}_{c_i})|$ **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

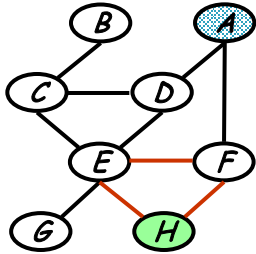
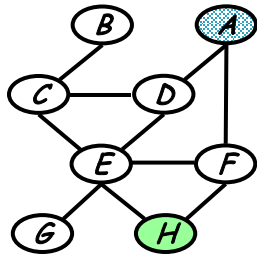
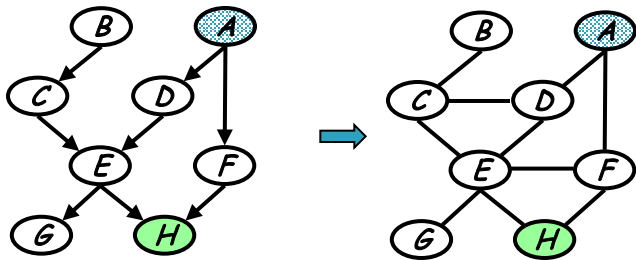
Complexity is **exponential** in number of variables in the intermediate factor

Understanding Variable Elimination

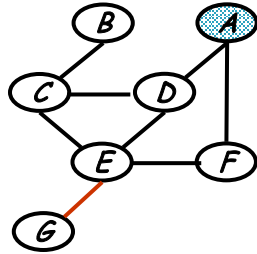
- ▶ A graph elimination algorithm



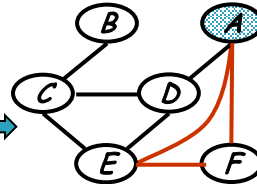
Elimination Cliques



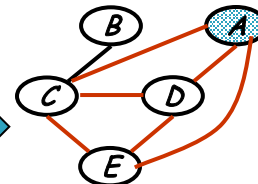
$m_h(e, f)$



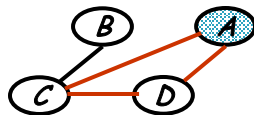
$m_g(e)$



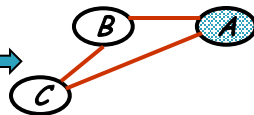
$m_f(e, a)$



$m_e(a, c, d)$



$m_d(a, c)$



$m_c(a, b)$

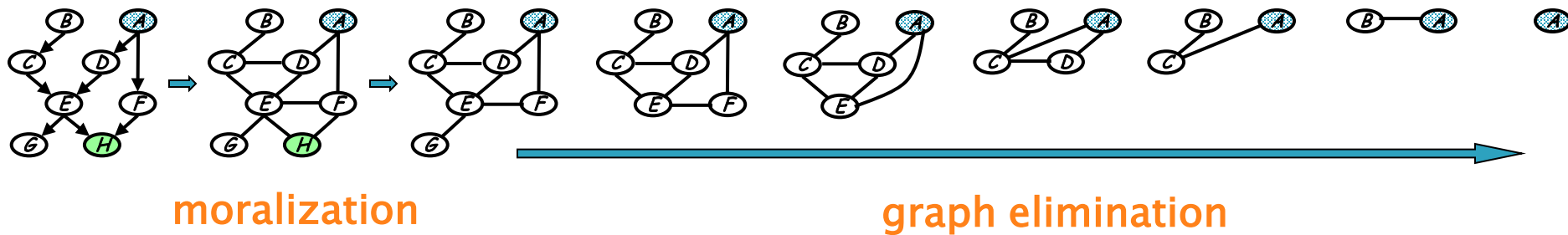


$m_b(a)$



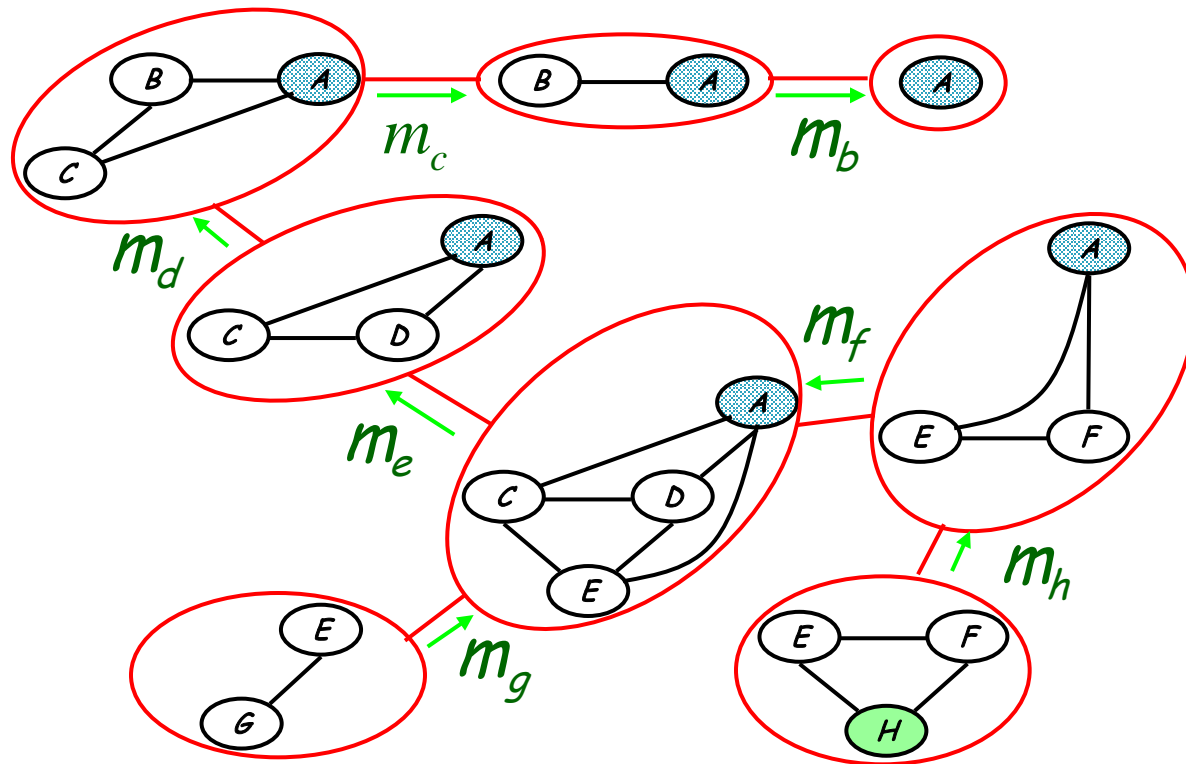
Understanding Variable Elimination

- ▶ A graph elimination algorithm



- ▶ Intermediate terms correspond to the **cliques** resulted from elimination
 - “good” elimination orderings lead to **small cliques** and hence reduce complexity (what will happen if we eliminate “e” first in the above graph?)
 - finding the optimum ordering is NP-hard, but for many graph optimum or near-optimum can often be heuristically found
- ▶ Applies to undirected GMs

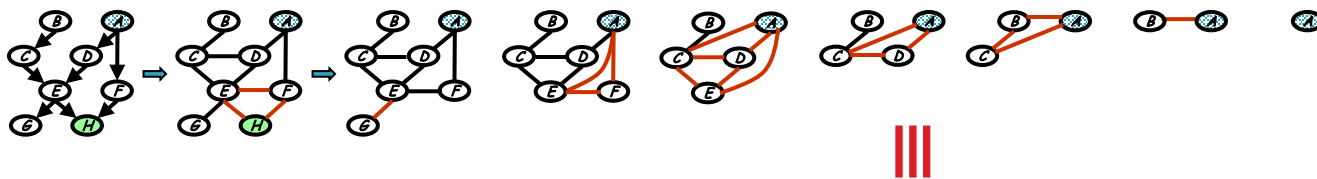
A clique tree



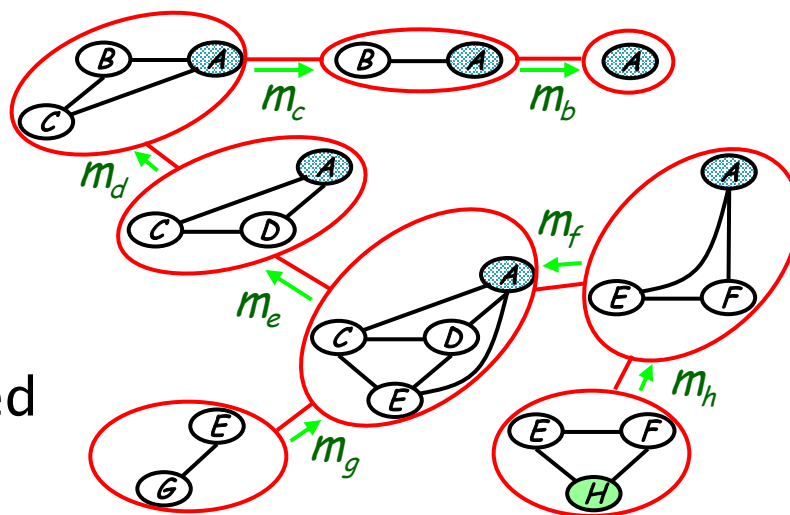
$$m_e(a, c, d) = \sum_e p(e | c, d) m_g(e) m_f(a, e)$$

From Elimination to Message Passing

- ▶ Our algorithm so far answers only one query (e.g., on one node), do we need to do a complete elimination for every such query?
- ▶ Elimination \equiv message passing on a **clique tree**



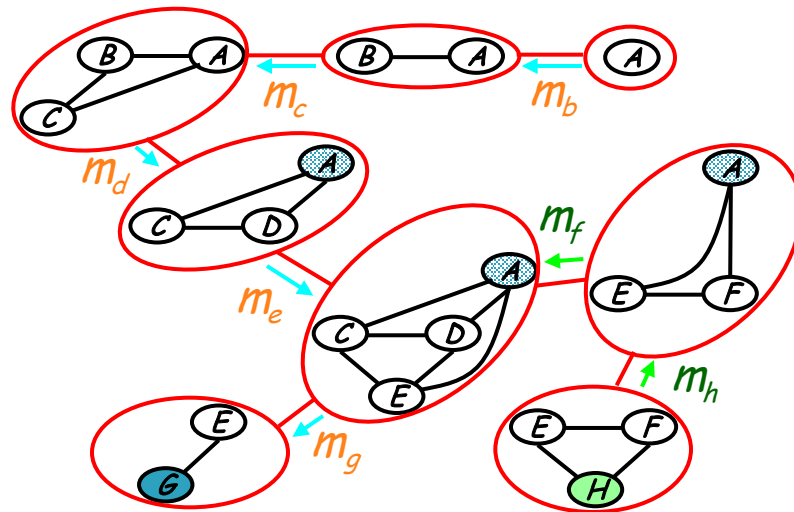
$$m_e(a, c, d) = \sum_e p(e | c, d) m_g(e) m_f(a, e)$$



- ▶ Messages can be reused

From Elimination to Message Passing

- ▶ Our algorithm so far answers only one query (e.g., on one node), do we need to do a complete elimination for every such query?
- ▶ Elimination \equiv message passing on a **clique tree**
 - **Another query ...**



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

A Sketch of the Junction Tree Algorithm

- ▶ **The algorithm**
 - Construction of junction trees --- a special **clique tree**
 - Propagation of probabilities --- a **message-passing protocol**
- ▶ Results in marginal probabilities of all cliques --- solves all queries in a single run
- ▶ A **generic** exact inference algorithm for any GM
- ▶ **Complexity**: exponential in the size of the maximal clique --- a good elimination order often leads to small maximal clique, and hence a good (i.e., thin) JT
- ▶ Many well-known algorithms are special cases of JT
 - Forward-backward, Kalman filter, Peeling, Sum-Product ...

Approaches to inference

- ▶ Exact inference algorithms
 - The elimination algorithm
 - The junction tree algorithms (not covered in detail)
- ▶ Approximate inference techniques
 - Stochastic simulation / sampling methods
 - Markov chain Monte Carlo methods
 - Variational algorithms (no covered)

Monte Carlo methods

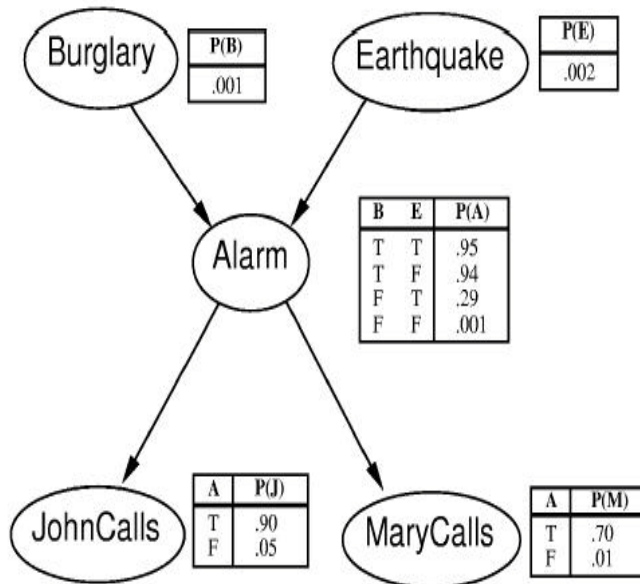
- ▶ Draw random samples from the desired distribution
- ▶ Yield a stochastic representation of a complex distribution
 - marginals and other expectations can be approximated using **sample-based averages**

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

- ▶ **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

- ▶ Sampling: Construct samples according to probabilities given in a BN.



E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
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, B_0 . Same for E_0 . $P(A|B_0, E_0) = \langle 0.001, 0.999 \rangle$ suppose it is false...

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

Example: naive sampling

- ▶ 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

Monte Carlo methods (cont.)

▶ Direct Sampling

- We have seen it.
- Very difficult to populate a high-dimensional state space

▶ Rejection Sampling

- Create samples like direct sampling, only count samples which is consistent with given evidence.

▶

▶ Markov chain Monte Carlo (MCMC)

Markov chain Monte Carlo

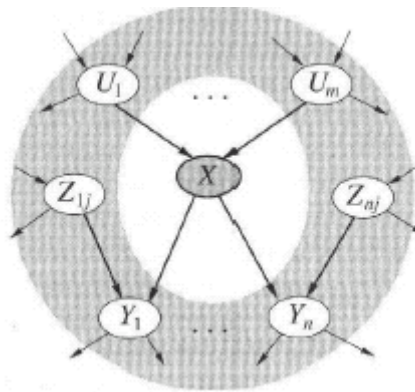
- ▶ Samples are obtained from a **Markov chain** (of sequentially evolving distributions) whose **stationary distribution** is the desired $p(x)$
- ▶ Gibbs sampling
 - we have variable set to $\mathbf{X} = \{x_1, x_2, x_3, \dots, x_N\}$
 - at each step one of the variables x_i is selected (at random or according to some fixed sequences)
 - the **conditional distribution** $p(x_i | \mathbf{x}_{-i})$ is computed
 - a value x_i is sampled from this distribution
 - the sample x_i replaces the previous of x_i in \mathbf{X} .

MCMC

▶ Markov-Blanket

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

$$\Rightarrow p(X_i | \mathbf{X}_{-i}) = p(X_i | \text{MB}(X_i))$$



▶ Gibbs sampling

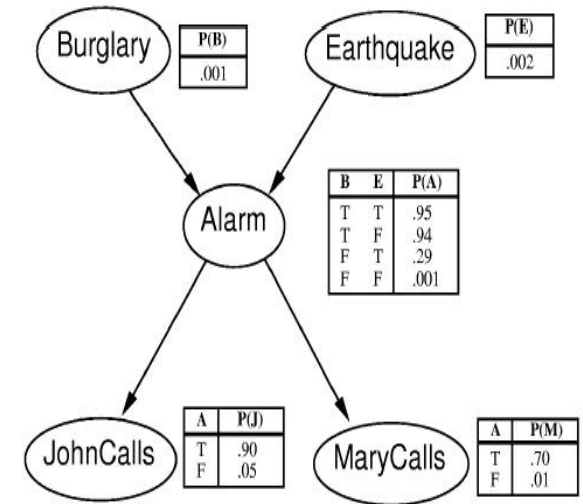
- Create a random sample. Every step, choose one variable and sample it by $P(X|\text{MB}(X))$ based on previous sample.

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

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

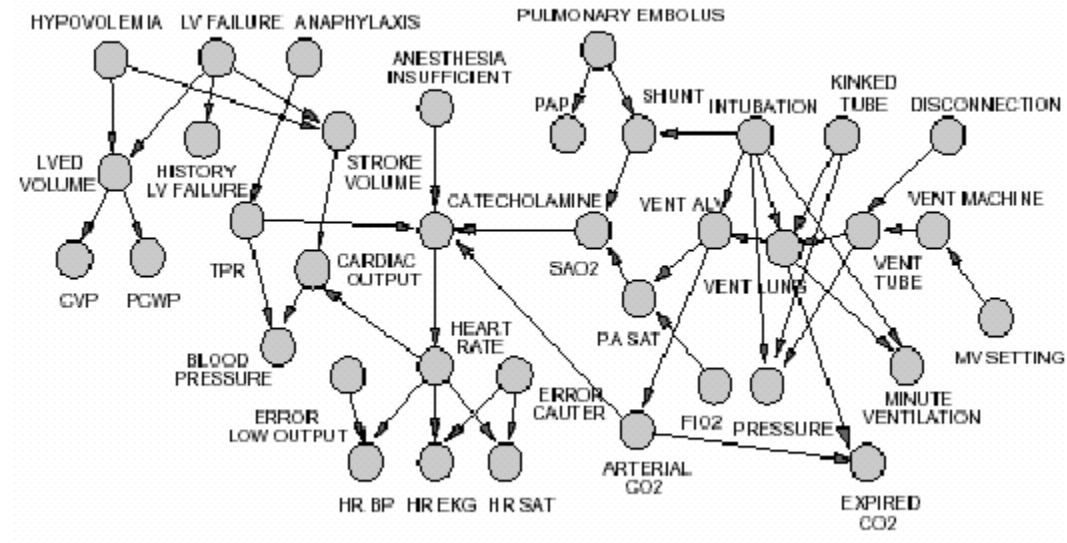
MCMC

- ▶ To calculate $P(J|B1,M1)$
- ▶ Choose $(B1,E0,A1,M1,J1)$ as a start
- ▶ Evidence is $B1, M1$, 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 $E \sim P(E|B1,A0)$
- ▶ ...



Complexity for Approximate Inference

- ▶ Approximate Inference will not reach the exact probability distribution in finite time, but only close to the value.
- ▶ Often much faster than exact inference when BN is **big** and **complex** enough. In MCMC, only consider $P(X|MB(X))$ but not the whole network.



Summary: inference

- ▶ Exact inference algorithms
 - The elimination algorithm
 - The junction tree algorithms (not covered in detail)
- ▶ Approximate inference techniques
 - Stochastic simulation / sampling methods
 - Markov chain Monte Carlo methods
 - Variational algorithms (no covered)
- ▶ Next Time....
 - Learning!!! Finally, ☺!