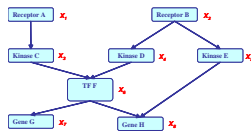


Machine Learning

10-701/15-781, Spring 2010

Bayesian Networks Learning and Inference



Eric Xing

Lecture 14, March 15, 2010
Reading: Chap. 8, C.B book

© Eric Xing @ CMU, 2006-2010

Recap of BN Representation



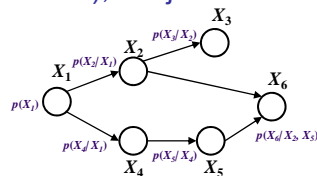
- Joint probability dist. on multiple variables:

$$P(X_1, X_2, X_3, X_4, X_5, X_6) = P(X_1)P(X_2 | X_1)P(X_3 | X_1, X_2)P(X_4 | X_1, X_2, X_3)P(X_5 | X_1, X_2, X_3, X_4)P(X_6 | X_1, X_2, X_3, X_4, X_5)$$

- If X_i 's are **independent**: $(P(X_i | \cdot) = P(X_i))$

$$P(X_1, X_2, X_3, X_4, X_5, X_6) = P(X_1)P(X_2)P(X_3)P(X_4)P(X_5)P(X_6) = \prod_i P(X_i)$$

- If X_i 's are **conditionally independent** (as described by a **GM**), the joint can be factored to simpler products, e.g.,



$$P(X_1, X_2, X_3, X_4, X_5, X_6) = P(X_1) P(X_2 | X_1) P(X_3 | X_2) P(X_4 | X_1) P(X_5 | X_4) P(X_6 | X_2, X_5)$$

© Eric Xing @ CMU, 2006-2010

2

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 ?
 - We use **learning** as a name for the process of obtaining point estimate of M .
 - But for *Bayesian*, they seek $p(M|D)$, which is actually an **inference** problem.
 - When not all variables are observable, even computing point estimate of M need to do **inference** to impute the *missing data*.

© Eric Xing @ CMU, 2006-2010

3

Learning BNs

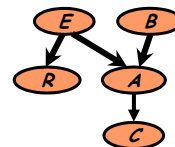


The goal:

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

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

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



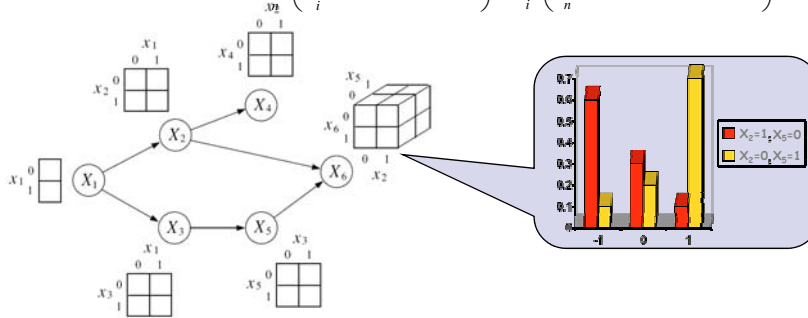
E	B	$P(A E, B)$	
e	b	0.9	0.1
e	\bar{b}	0.2	0.8
\bar{e}	b	0.9	0.1
\bar{e}	\bar{b}	0.01	0.99

© Eric Xing @ CMU, 2006-2010

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:

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



© Eric Xing @ CMU, 2006-2010

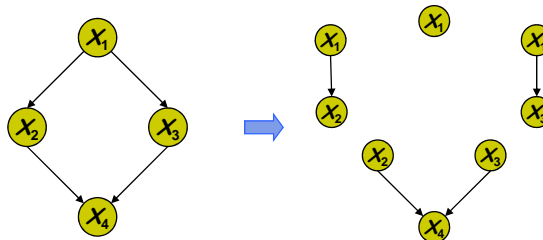
5

Example: decomposable likelihood of a directed model

- Consider the distribution defined by the directed acyclic GM:

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

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



© Eric Xing @ CMU, 2006-2010

6

E.g.: MLE for BNs with tabular CPDs



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

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

- Note that in case of multiple parents, \mathbf{x}_{π_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:

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



© Eric Xing @ CMU, 2006-2010

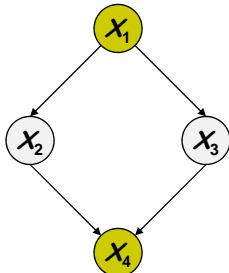
7

What if some nodes are not observed?



- Consider the distribution defined by the directed acyclic GM:

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



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

© Eric Xing @ CMU, 2006-2010

8

Probabilistic Inference



- **Computing statistical queries regarding the network, e.g.:**
 - Is node X independent on node Y given nodes Z, W ?
 - What is the probability of $X=\text{true}$ if ($Y=\text{false}$ and $Z=\text{true}$)?
 - What is the joint distribution of (X, Y) if $Z=\text{false}$?
 - What is the likelihood of some full assignment?
 - What is the most likely assignment of values to all or a subset the nodes of the network?
- **General purpose algorithms exist to fully automate such computation**
 - Computational cost depends on the topology of the network
 - **Exact inference:**
 - The junction tree algorithm
 - **Approximate inference;**
 - Loopy belief propagation, variational inference, Monte Carlo sampling

Inferential Query 1: Likelihood



- Most of the queries one may ask involve **evidence**
 - Evidence \mathbf{x}_v is an assignment of values to a set \mathbf{X}_v of nodes in the GM over variable set $\mathbf{X}=\{X_1, X_2, \dots, X_n\}$
 - Without loss of generality $\mathbf{X}_v=\{X_{k+1}, \dots, X_n\}$,
 - Write $\mathbf{X}_H=\mathbf{X}\setminus\mathbf{X}_v$ as the set of hidden variables, \mathbf{X}_H can be \emptyset or \mathbf{X}
- Simplest query: compute probability of evidence

$$P(\mathbf{x}_v) = \sum_{\mathbf{X}_H} P(\mathbf{X}_H, \mathbf{x}_v) = \sum_{x_1} \dots \sum_{x_k} P(x_1, \dots, x_k, \mathbf{x}_v)$$

- this is often referred to as computing the **likelihood** of \mathbf{x}_v

Inferential Query 2: Conditional Probability



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

$$P(\mathbf{X}_H | \mathbf{X}_V = \mathbf{x}_V) = \frac{P(\mathbf{X}_H, \mathbf{X}_V)}{P(\mathbf{X}_V)} = \frac{P(\mathbf{X}_H, \mathbf{X}_V)}{\sum_{\mathbf{x}_H} P(\mathbf{X}_H = \mathbf{x}_H, \mathbf{X}_V)}$$

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

$$P(\mathbf{Y} | \mathbf{x}_V) = \sum_{\mathbf{z}} P(\mathbf{Y}, \mathbf{Z} = \mathbf{z} | \mathbf{x}_V)$$

- the process of summing out the "don't care" variables \mathbf{z} is called **marginalization**, and the resulting $P(\mathbf{Y} | \mathbf{x}_V)$ is called a **marginal prob.**

© Eric Xing @ CMU, 2006-2010

11

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 form all parts of the network

© Eric Xing @ CMU, 2006-2010

12

Inferential 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 \mathbf{x}_v , and ignoring (the values of) other variables \mathbf{Z} :

$$\mathbf{Y}^* | \mathbf{x}_v = \arg \max_{\mathbf{y}} P(\mathbf{Y} | \mathbf{x}_v) = \arg \max_{\mathbf{y}} \sum_{\mathbf{z}} P(\mathbf{Y}, \mathbf{Z} = \mathbf{z} | \mathbf{x}_v)$$

- this is the **maximum a posteriori** configuration of \mathbf{Y} .

Complexity of Inference



Thm:

Computing $P(X_H = x_H | \mathbf{x}_v)$ in an arbitrary BN 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 BNs
 - For particular families of BNs, we can have provably efficient procedures

Approaches to inference



- Exact inference algorithms
 - The elimination algorithm ✓
 - The junction tree algorithms ✓ (but will not cover in detail here)
- Approximate inference techniques
 - Stochastic simulation / sampling methods ✓
 - Markov chain Monte Carlo methods ✓
 - Variational algorithms (will be covered in advanced ML courses)

Marginalization and Elimination



- A signal transduction pathway:

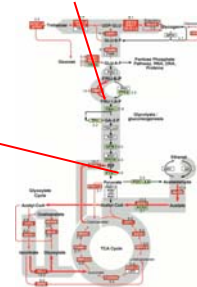


What is the likelihood that protein E is active?

- Query: $P(e)$

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

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



- By chain decomposition, we get

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

Elimination on Chains



- Rearranging terms ...

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

© Eric Xing @ CMU, 2006-2010

17

Elimination on Chains



- Now we can perform innermost summation

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

- This summation "eliminates" one variable from our summation argument at a "local cost".

© Eric Xing @ CMU, 2006-2010

18

Elimination in Chains



- Rearranging and then summing again, we get

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

© Eric Xing @ CMU, 2006-2010

19

Elimination in Chains



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

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

- Complexity:
 - Each step costs $O(|Val(X_i)| * |Val(X_{i+1})|)$ operations: $O(nk^2)$
 - Compare to naïve evaluation that sums over joint values of $n-1$ variables $O(k^n)$

© Eric Xing @ CMU, 2006-2010

20

Inference on General BN via Variable Elimination



General idea:

- Write query in the form

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

- this suggests an "elimination order" of latent 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

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

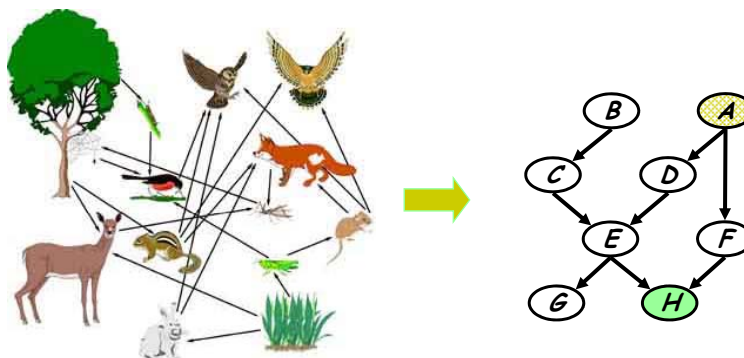
© Eric Xing @ CMU, 2006-2010

21

A more complex network



A food web



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

© Eric Xing @ CMU, 2006-2010

22

Example: Variable Elimination

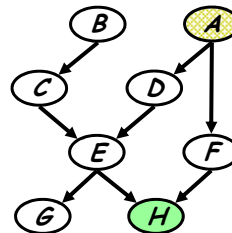


- Query: $P(A | h)$
 - Need to eliminate: B, C, D, E, F, G, H

- 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



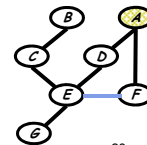
- Step 1:

- **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:

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



© Eric Xing @ CMU, 2006-2010

23

Example: Variable Elimination



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

- Initial factors:

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

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

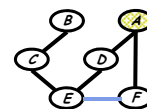
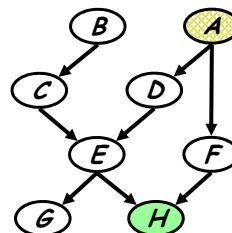
- Step 2: Eliminate G

- compute

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

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

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



© Eric Xing @ CMU, 2006-2010

24

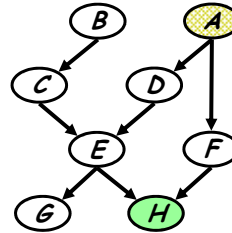
Example: Variable Elimination



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

- Initial factors:

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

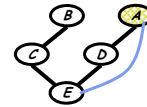


- Step 3: Eliminate F

- compute

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

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



© Eric Xing @ CMU, 2006-2010

25

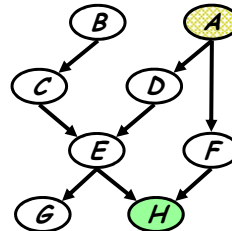
Example: Variable Elimination



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

- Initial factors:

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

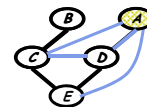


- Step 4: Eliminate E

- compute

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

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



© Eric Xing @ CMU, 2006-2010

26

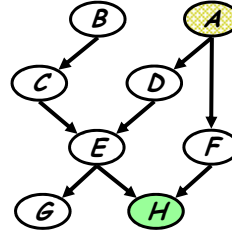
Example: Variable Elimination



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

- Initial factors:

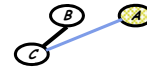
$$\begin{aligned}
 &P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)P(g|e)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)P(e|c,d)P(f|a)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)P(e|c,d)m_f(a,e) \\
 \Rightarrow &P(a)P(b)P(c|b)P(d|a)m_e(a,c,d)
 \end{aligned}$$



- Step 5: Eliminate D

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

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



© Eric Xing @ CMU, 2006-2010

27

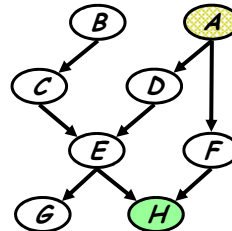
Example: Variable Elimination



- 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 $m_c(a,b) = \sum_c p(c|b)m_d(a,c)$

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



© Eric Xing @ CMU, 2006-2010

28

Example: Variable Elimination



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

- Initial factors:

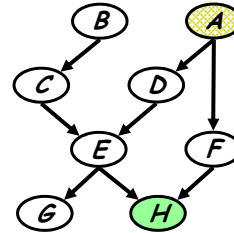
$$\begin{aligned}
 &P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)m_h(e,f) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)P(e|c,d)m_f(a,e) \\
 \Rightarrow &P(a)P(b)P(c|d)P(d|a)m_e(a,c,d) \\
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

- Step 7: Eliminate B

- compute

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

$$\Rightarrow P(a)m_b(a)$$



(A)

© Eric Xing @ CMU, 2006-2010

29

Example: Variable Elimination



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

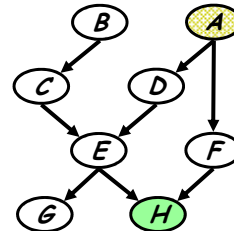
- Initial factors:

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

- Step 8: Wrap-up

$$p(a, \tilde{h}) = p(a)m_b(a), \quad p(\tilde{h}) = \sum_a p(a)m_b(a)$$

$$\Rightarrow P(a | \tilde{h}) = \frac{p(a)m_b(a)}{\sum_a p(a)m_b(a)}$$



© Eric Xing @ CMU, 2006-2010

30

Complexity of variable elimination



- Suppose in one elimination step we compute

$$m_x(y_1, \dots, y_k) = \sum_x m'_x(x, y_1, \dots, y_k)$$

$$m'_x(x, y_1, \dots, y_k) = \prod_{i=1}^k m_i(x, y_{c_i})$$

This requires

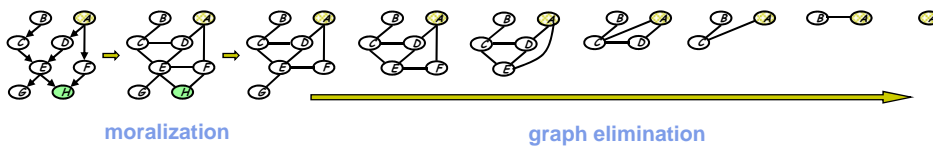
- $k \cdot |\text{Val}(X)| \cdot \prod_i |\text{Val}(Y_{c_i})|$ **multiplications**
 - For each value of x, y_1, \dots, y_k , we do k multiplications
- $|\text{Val}(X)| \cdot \prod_i |\text{Val}(Y_{c_i})|$ **additions**
 - For each value of y_1, \dots, y_k , we do $|\text{Val}(X)|$ additions

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

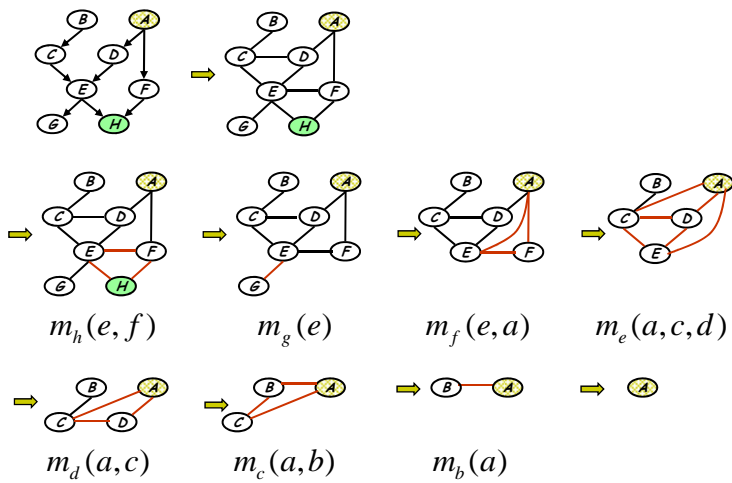
Understanding Variable Elimination



- A graph elimination algorithm



Elimination Cliques



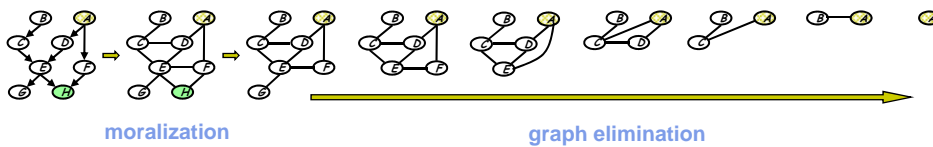
© Eric Xing @ CMU, 2006-2010

33

Understanding Variable Elimination



- A graph elimination algorithm

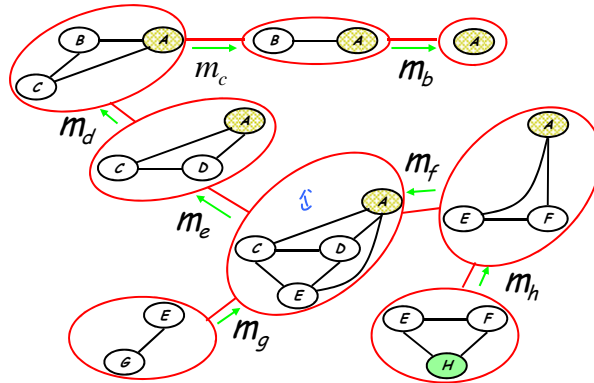


- Intermediate terms correspond to the **cliques** resulted from elimination
 - "good" elimination orderings lead to **small cliques** and hence reduce complexity (what will happen if we eliminate "e" first in the above graph?)
 - finding the optimum ordering is NP-hard, but for many graph optimum or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

© Eric Xing @ CMU, 2006-2010

34

A clique tree



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

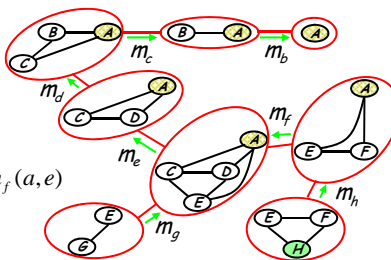
© Eric Xing @ CMU, 2006-2010

35

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

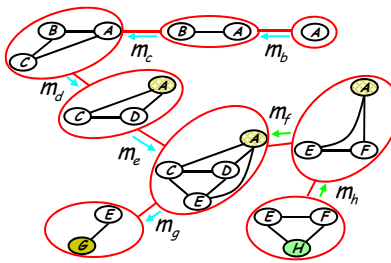
© Eric Xing @ CMU, 2006-2010

36

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

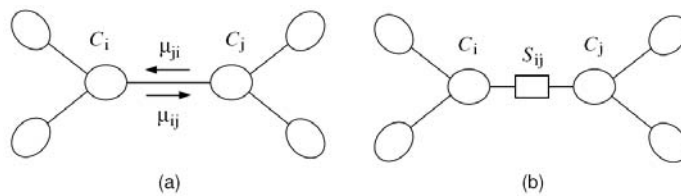
© Eric Xing @ CMU, 2006-2010

37

The Junction Tree Algorithm



- Shafer-Shenoy algorithm



- Message from clique i to clique j :

$$\mu_{i \rightarrow j} = \sum_{C_i \setminus S_{ij}} \psi_{C_i} \prod_{k \neq j} \mu_{k \rightarrow i}(S_{ki})$$

- Clique marginal

$$p(C_i) \propto \psi_{C_i} \prod_k \mu_{k \rightarrow i}(S_{ki})$$

© Eric Xing @ CMU, 2006-2010

38

A Sketch of the Junction Tree Algorithm

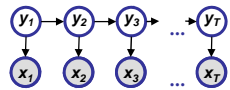


- The algorithm
 - Construction of junction trees --- a special **clique tree**
 - Propagation of probabilities --- a **message-passing protocol**
- Results in marginal probabilities of all cliques --- solves all queries in a single run
- A **generic** exact inference algorithm for any GM
- **Complexity**: 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 ...

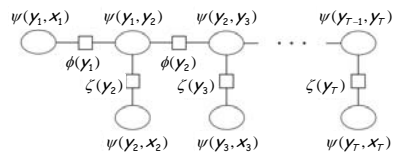
The Junction tree algorithm for HMM



- A junction tree for the HMM



⇒



- Rightward pass

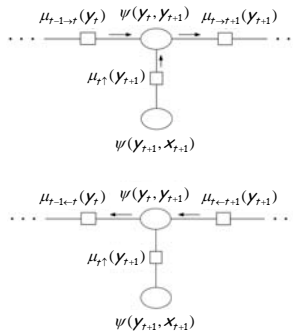
$$\begin{aligned} \mu_{t \rightarrow t+1}(y_{t+1}) &= \sum_{y_t} \psi(y_t, y_{t+1}) \mu_{t-1 \rightarrow t}(y_t) \mu_{t \uparrow}(y_{t+1}) \\ &= \sum_{y_t} p(y_{t+1} | y_t) \mu_{t-1 \rightarrow t}(y_t) p(x_{t+1} | y_{t+1}) \\ &= p(x_{t+1} | y_{t+1}) \sum_{y_t} a_{y_t, y_{t+1}} \mu_{t-1 \rightarrow t}(y_t) \end{aligned}$$

- This is exactly the **forward algorithm!**

- Leftward pass ...

$$\begin{aligned} \mu_{t-1 \leftarrow t}(y_t) &= \sum_{y_{t+1}} \psi(y_t, y_{t+1}) \mu_{t \leftarrow t+1}(y_{t+1}) \mu_{t \uparrow}(y_{t+1}) \\ &= \sum_{y_{t+1}} p(y_{t+1} | y_t) \mu_{t \leftarrow t+1}(y_{t+1}) p(x_{t+1} | y_{t+1}) \end{aligned}$$

- This is exactly the **backward algorithm!**



Summary



- Represent dependency structure with a directed acyclic graph
 - Node \leftrightarrow random variable
 - Edges encode dependencies
 - Absence of edge \rightarrow conditional independence
 - Plate representation
 - A BN is a database of prob. Independence statement on variables
- The factorization theorem of the joint probability
 - Local specification \rightarrow globally consistent distribution
 - Local representation for exponentially complex state-space
- Support efficient inference and learning

