

Readings:

K&F: 8.1, 8.2, 8.3, 8.7.1

K&F: 9.1, 9.2, 9.3, 9.4

Variable Elimination 2

Clique Trees

Graphical Models – 10708

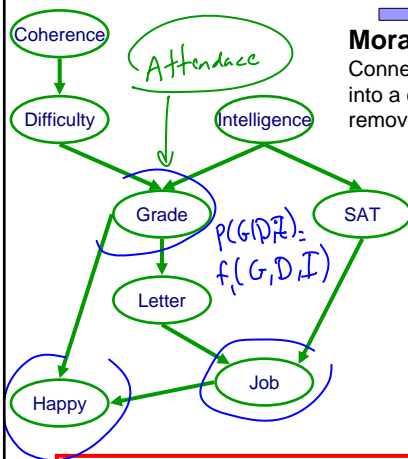
Carlos Guestrin

Carnegie Mellon University

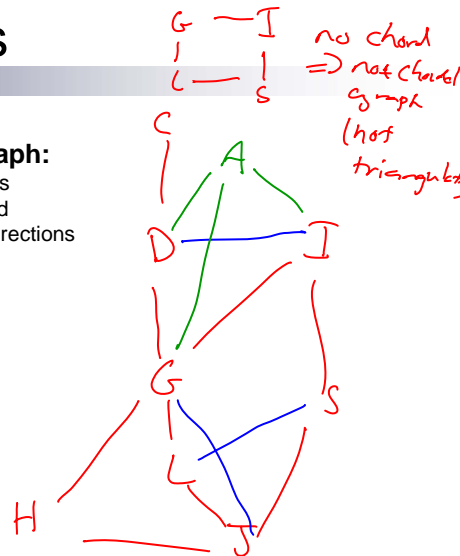
October 13th, 2006

1

Complexity of variable elimination – Graphs with loops



Moralize graph:
Connect parents
into a clique and
remove edge directions



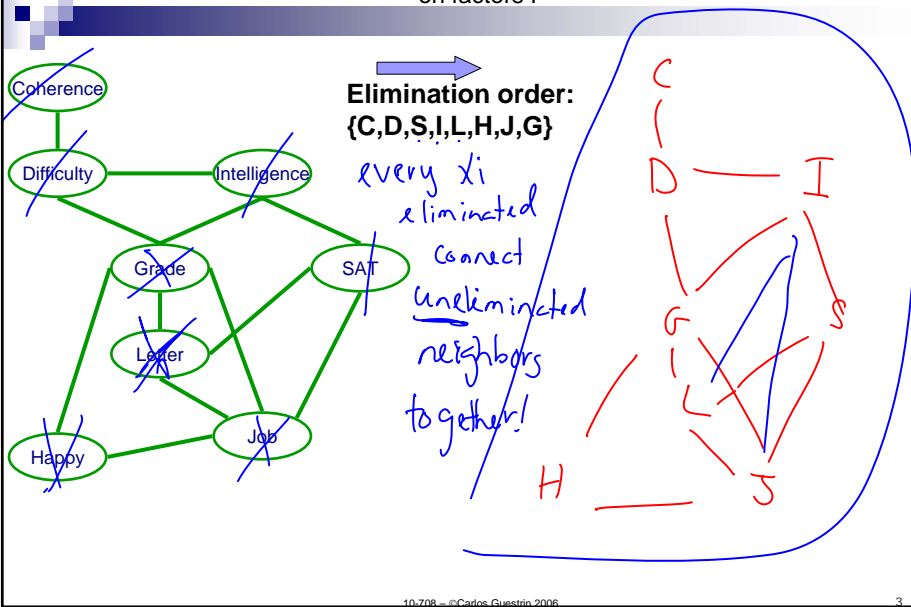
Connect nodes that appear together in an initial factor

10708 – ©Carlos Guestrin 2006

2

Induced graph

The **induced graph** $I_{F \prec}$ for elimination order \prec has an edge $X_i - X_j$ if X_i and X_j appear together in a factor generated by VE for elimination order \prec on factors F

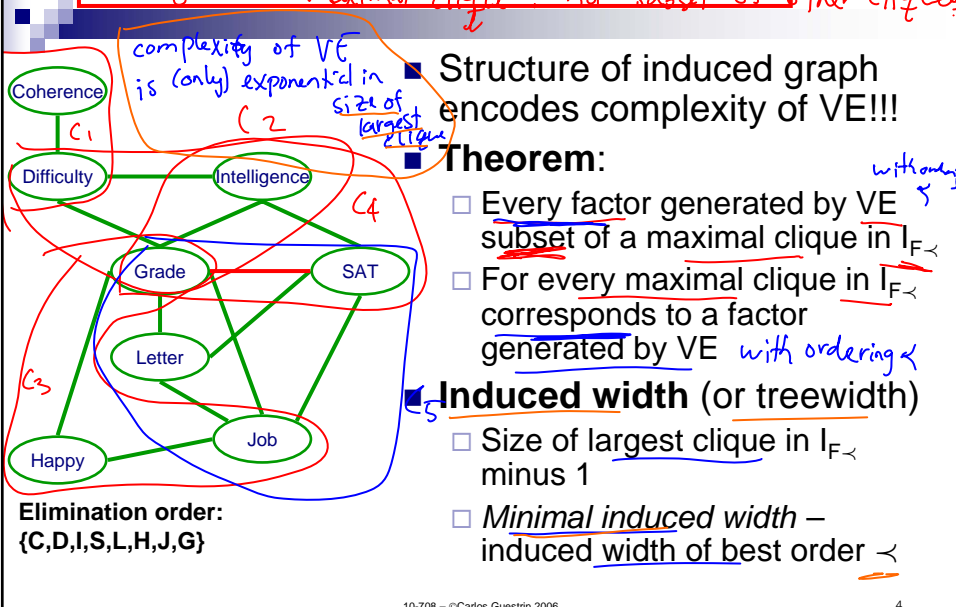


10.708 - ©Carlos Guestrin 2006

3

Induced graph and complexity of VE

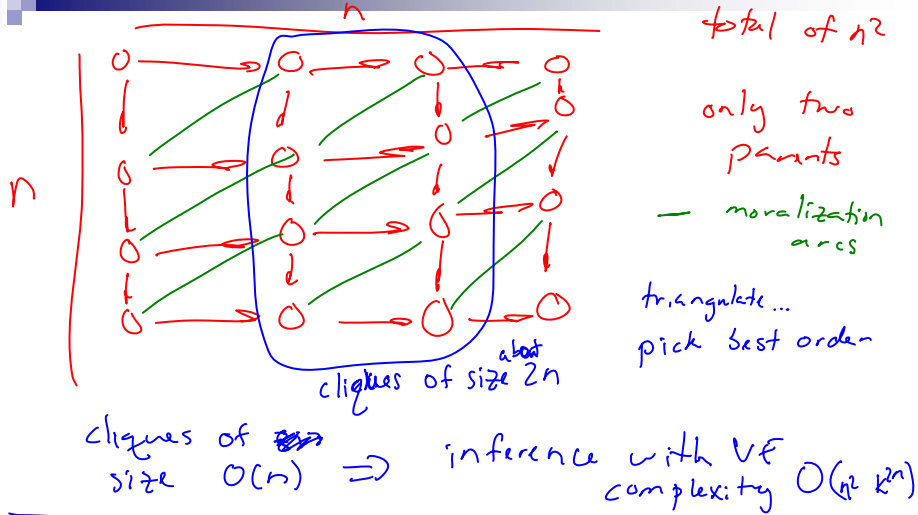
Read complexity from cliques in induced graph



10.708 - ©Carlos Guestrin 2006

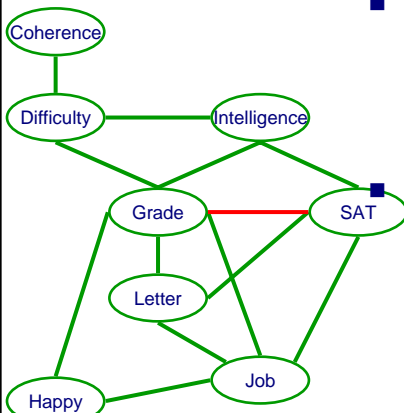
4

Example: Large induced-width with small number of parents



Compact representation \nRightarrow Easy inference ☹️

Finding optimal elimination order



Elimination order:
{C,D,I,S,L,H,J,G}

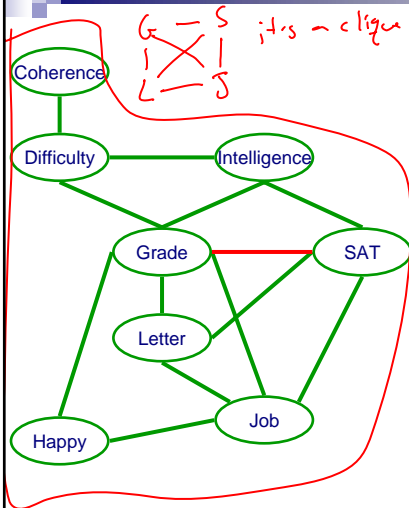
■ **Theorem:** Finding best elimination order is NP-complete:

- Decision problem: Given a graph, determine if there exists an elimination order that achieves induced width $\leq K$

■ **Interpretation:**

- Hardness of finding elimination order in addition to hardness of inference
- Actually, can find elimination order in time exponential in size of largest clique – same complexity as inference

Induced graphs and chordal graphs



Chordal graph:

- Every cycle $X_1 - X_2 - \dots - X_k - X_1$ with $k \geq 3$ has a chord
- Edge $X_i - X_j$ for non-consecutive i & j

Theorem:

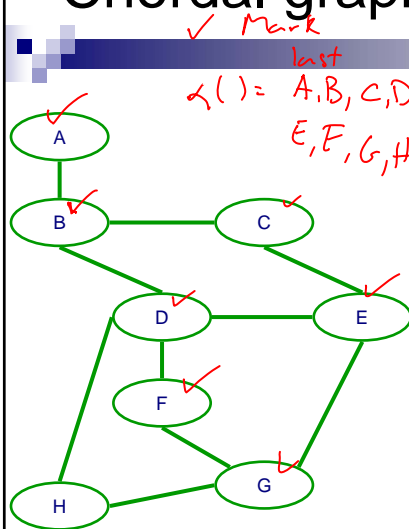
- Every induced graph is chordal

- “Optimal” elimination order easily obtained for chordal graph

10.708 - ©Carlos Guestrin 2006

7

Chordal graphs and triangulation



Triangulation: turning graph into chordal graph

Max Cardinality Search:

- Simple heuristic
- Initialize unobserved nodes X as unmarked
- For $k = |X|$ to 1
 - $X \leftarrow$ unmarked var with most marked neighbors
 - $\prec(X) \leftarrow k$
 - Mark X

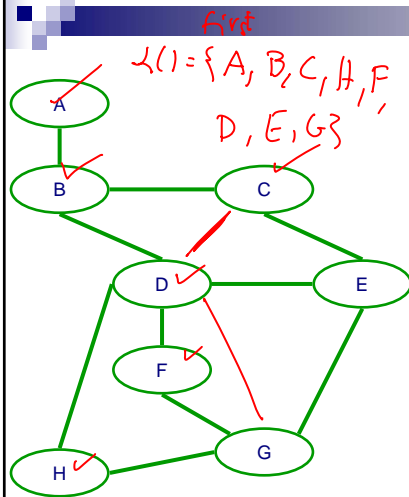
Theorem: Obtains optimal order for chordal graphs

- Often, not so good in other graphs!

10.708 - ©Carlos Guestrin 2006

8

Minimum fill/size/weight heuristics



- Many more effective heuristics
 - see reading
- **Min (weighted) fill heuristic**
 - Often very effective
- Initialize unobserved nodes **X** as unmarked
- For $k = 1$ to $|X|$
 - $X \leftarrow$ unmarked var whose elimination adds fewest edges
 - $\prec(X) \leftarrow k$
 - Mark X
 - Add fill edges introduced by eliminating X
- **Weighted version:**
 - Consider size of factor rather than number of edges

10.708 - ©Carlos Guestrin 2006

9

Choosing an elimination order

- Choosing best order is NP-complete
 - Reduction from MAX-Clique
- Many good heuristics (some with guarantees)
- Ultimately, can't beat NP-hardness of inference
 - Even optimal order can lead to exponential variable elimination computation
- In practice
 - Variable elimination often very effective
 - Many (many many) approximate inference approaches available when variable elimination too expensive
 - Most approximate inference approaches build on ideas from variable elimination

10.708 - ©Carlos Guestrin 2006

10

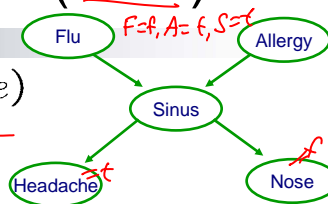
Announcements

- Recitation on advanced topic:
 - Carlos on Context-Specific Independence
 - On Monday Oct 16, 5:30-7:00pm in Wean Hall 4615A

○ HW3 out later today

Most likely explanation (MLE)

- Query: $\text{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n | e)$



- Using defn of conditional probs:

$$\text{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n | e) = \text{argmax}_{x_1, \dots, x_n} \frac{P(x_1, \dots, x_n, e)}{P(e)}$$

- Normalization irrelevant:

$$\text{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n | e) = \text{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n, e)$$

Max-marginalization



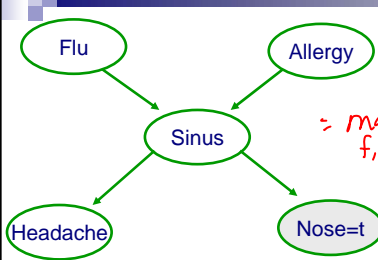
$$\max_{f,s} P(f) \cdot P(s|f) \cdot P(N=t|s)$$

$$\max \left[\begin{aligned} &P(f=t) \cdot P(s=t|f=t) \cdot P(N=t|s=t); \\ &P(f=t) \cdot P(s=t|f=t) \cdot P(N=t|s=f); \\ &P(f=t) \cdot P(s=f|f=t) \cdot P(N=t|s=f); \\ &P(f=f) \cdot P(s=f|f=f) \cdot P(N=t|s=f); \\ &P(f=f) \cdot P(s=f|f=f) \cdot P(N=t|s=t); \end{aligned} \right]$$

10.708 - ©Carlos Guestrin 2006

13

Example of variable elimination for MLE - Forward pass



$$\max_{f,a,s,h} P(f) \cdot P(a) \cdot P(s|f,a) \cdot P(h|s) \cdot P(N=t|s)$$

$$= \max_{f,a,s} P(f) \cdot P(a) \cdot P(s|f,a) \cdot P(N=t|s) \cdot \max_h P(h|s)$$

$$g_1(s)$$

$$= \max_{f,a} P(f) \cdot P(a) \cdot \max_s P(s|f,a) \cdot P(N=t|s) \cdot g_1(s)$$

$$g_2(f,a)$$

$$= \max_f P(f) \cdot \max_a P(a) \cdot g_2(f,a)$$

$$g_3(f)$$

$$g_4 = \max_f P(f) \cdot g_3(f)$$

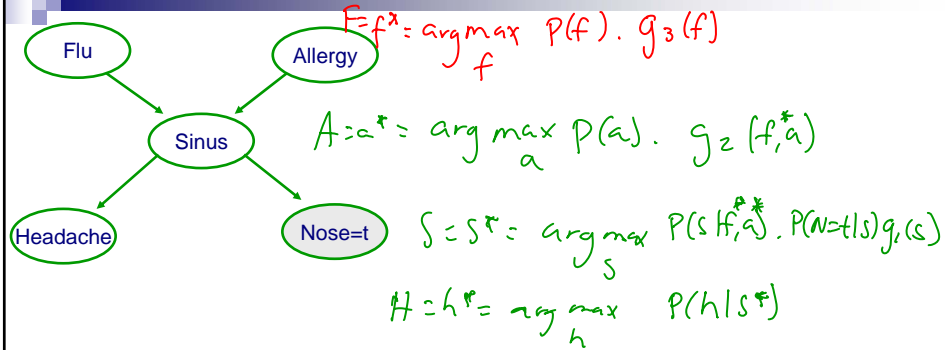
Interpretation:

- $g_1(s)$ = highest prob h achieves for each value of s
- $g_2(f,a)$ = highest prob. s & h | $N=t$ for each value of f, a
- g_4 = prob. MLE $N=t$

10.708 - ©Carlos Guestrin 2006

14

Example of variable elimination for MLE – Backward pass



10.708 – ©Carlos Guestrin 2006

15

MLE Variable elimination algorithm – Forward pass

use max instead of sum (no pruning)

- Given a BN and a MLE query $\max_{x_1, \dots, x_n} P(x_1, \dots, x_n, \mathbf{e})$
- Instantiate evidence $\mathbf{E} = \mathbf{e}$
- Choose an ordering on variables, e.g., X_1, \dots, X_n
- For $i = 1$ to n , If $X_i \notin \mathbf{E}$
 - Collect factors f_1, \dots, f_k that include X_i
 - Generate a new factor by eliminating X_i from these factors

$$g = \max_{x_i} \prod_{j=1}^k f_j$$

- Variable X_i has been eliminated!
- cache g*

10.708 – ©Carlos Guestrin 2006

16

MLE Variable elimination algorithm – Backward pass

- $\{x_1^*, \dots, x_n^*\}$ will store maximizing assignment
- For $i = n$ to 1, If $X_i \notin \mathbf{E}$
 - Take factors f_1, \dots, f_k used when X_i was eliminated
 - Instantiate f_1, \dots, f_k , with $\{x_{i+1}^*, \dots, x_n^*\}$
 - Now each f_j depends only on X_i
 - Generate maximizing assignment for X_i :

$$x_i^* \in \operatorname{argmax}_{x_i} \prod_{j=1}^k f_j$$

10.708 – ©Carlos Guestrin 2006

17

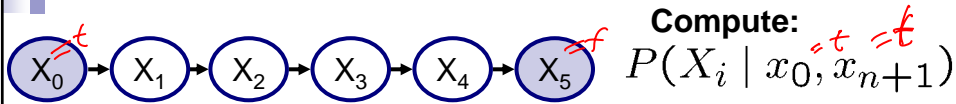
What you need to know about VE

- Variable elimination algorithm
 - Eliminate a variable:
 - Combine factors that include this var into single factor
 - Marginalize var from new factor
 - Cliques in induced graph correspond to factors generated by algorithm
 - Efficient algorithm (“only” exponential in induced-width, not number of variables)
 - If you hear: “Exact inference only efficient in tree graphical models”
 - You say: “No!!! Any graph with low induced width”
 - And then you say: “And even some with very large induced-width” (special recitation)
- Elimination order is important!
 - NP-complete problem
 - Many good heuristics
- Variable elimination for MLE
 - Only difference between probabilistic inference and MLE is “sum” versus “max”

10.708 – ©Carlos Guestrin 2006

18

What if I want to compute $P(X_i | x_0, x_{n+1})$ for each i ?



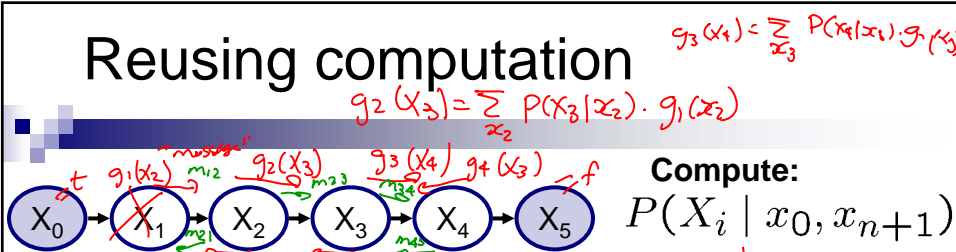
Variable elimination for each i ? e.g., $x_1 \dots x_{i-1}, x_{i+1} \dots x_n$

eliminate x_1 $g_1(x_2) = \sum_{x_1} P(x_0) \cdot P(x_1 | x_0) \cdot P(x_2 | x_1)$
 complexity of $P(x_i | x_0, x_{n+1}) = O(n)$

Variable elimination for every i , what's the complexity?

naive $O(n^2)$ | run VE n times:

Reusing computation



$P(x_3 | x_0, x_5)$ - compute $g_1(x_2), g_2(x_3), g_3(x_4)$ (doesn't need $g_3(x_4)$)
 $P(x_4 | x_0, x_5)$ - compute $g_1(x_2), g_2(x_3), g_3(x_4)$ (same g_1)

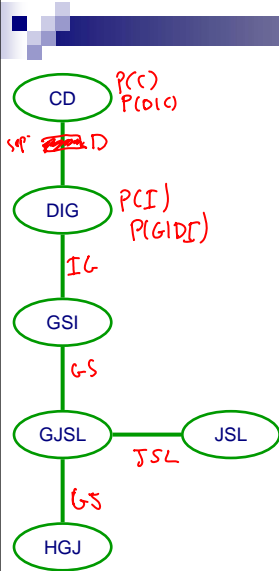
done! $P(x_4 | x_0, x_5) \propto g_3(x_4) \cdot P(x_5 = f)$

need to eliminate x_4 : done!

$$g_4(x_3) = \sum_{x_4} P(x_5 = f | x_4) \cdot P(x_4 | x_3)$$

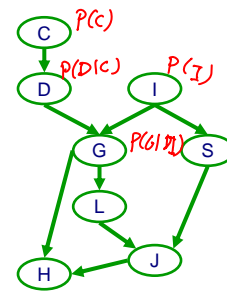
compute each message once!! , so two passes ($O(n)$) gives you all probs.

Cluster graph

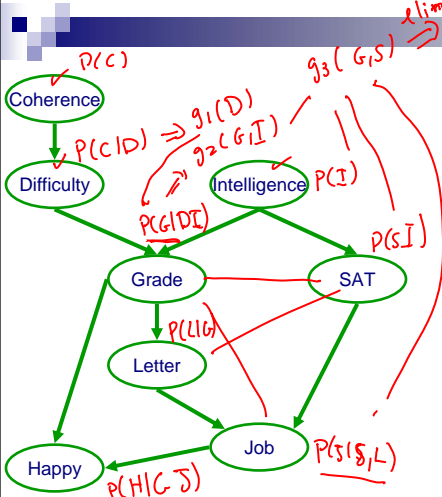


Cluster graph: For set of factors F

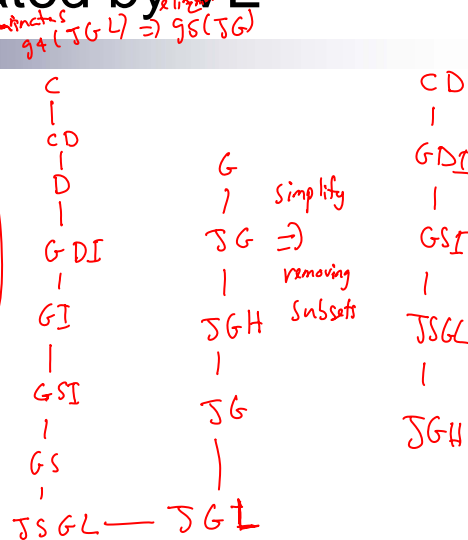
- Undirected graph
- Each node i associated with a cluster C_i
- Family preserving: for each factor $f_j \in F$, \exists node i such that $\text{scope}[f_j] \subseteq C_i$
- Each edge $i - j$ is associated with a separator $S_{ij} = C_i \cap C_j$



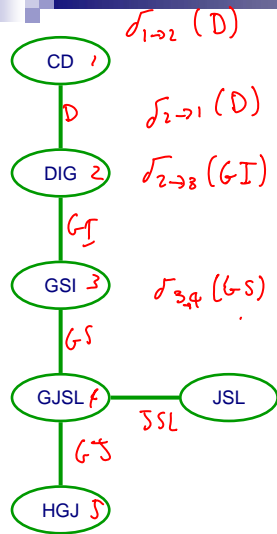
Factors generated by VE



Elimination order:
 $\{C, D, I, S, L, H, J, G\}$



Cluster graph for VE



VE generates cluster tree!

- One clique for each factor used/generated
- Edge $i - j$, if f_i used to generate f_j
- “Message” from i to j generated when marginalizing a variable from f_i
- Tree because factors only used once

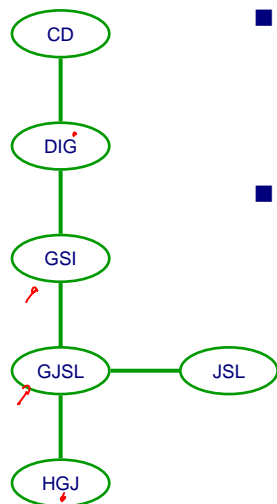
Proposition:

- “Message” δ_{ij} from i to j
- Scope $[\delta_{ij}] \subseteq \mathbf{S}_{ij}$

10.708 – ©Carlos Guestrin 2006

23

Running intersection property



Running intersection property (RIP)

- Cluster tree satisfies RIP if whenever $X \in \mathbf{C}_i$ and $X \in \mathbf{C}_j$ then X is in every cluster in the (unique) path from \mathbf{C}_i to \mathbf{C}_j

Theorem:

- Cluster tree generated by VE satisfies RIP

10.708 – ©Carlos Guestrin 2006

24

Constructing a clique tree from VE

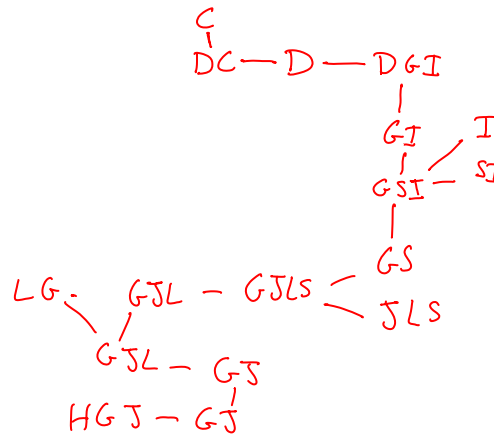
- Select elimination order

<

- Connect factors that would be generated if you run VE with order <

- Simplify!

- Eliminate factor that is subset of neighbor



10.708 - ©Carlos Guestrin 2006

25

Find clique tree from chordal graph

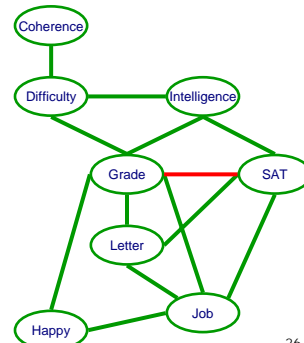
- Triangulate moralized graph to obtain chordal graph

- Find maximal cliques

- NP-complete in general
- Easy for chordal graphs
- Max-cardinality search

- Maximum spanning tree finds clique tree satisfying RIP!!!

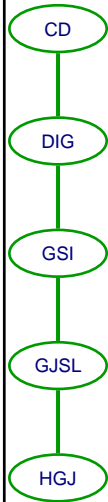
- Generate weighted graph over cliques
- Edge weights (i,j) is separator size - $|C_i \cap C_j|$



10.708 - ©Carlos Guestrin 2006

26

Clique tree & Independencies



■ Clique tree (or Junction tree)

- A cluster tree that satisfies the RIP

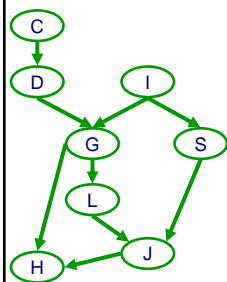
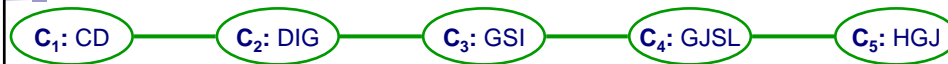
■ Theorem:

- Given some BN with structure G and factors F
- For a clique tree T for F consider $C_i - C_j$ with separator S_{ij} :
 - X – any set of vars in C_i side of the tree
 - Y – any set of vars in C_j side of the tree
- Then, $(X \perp Y \mid S_{ij})$ in BN
- Furthermore, $I(T) \subseteq I(G)$

10.708 – ©Carlos Guestrin 2006

27

Variable elimination in a clique tree 1



■ Clique tree for a BN

- Each CPT assigned to a clique
- Initial potential $\pi_0(C_i)$ is product of CPTs

10.708 – ©Carlos Guestrin 2006

28

Variable elimination in a clique tree 2



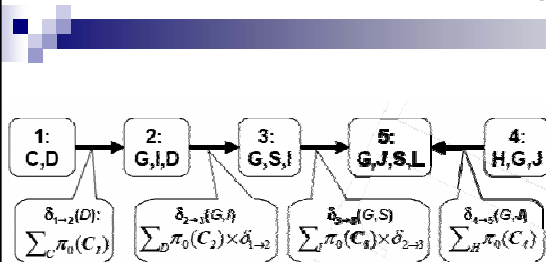
■ VE in clique tree to compute $P(X_i)$

- Pick a root (any node containing X_i)
- Send messages recursively from leaves to root
 - Multiply incoming messages with initial potential
 - Marginalize vars that are not in separator
- Clique *ready* if received messages from all neighbors

10.708 - ©Carlos Guestrin 2006

29

Belief from message



■ Theorem: When clique C_i is ready

- Received messages from all neighbors
- Belief $\pi_i(C_i)$ is product of initial factor with messages:

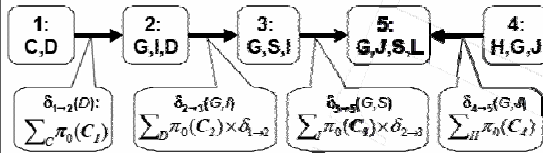
10.708 - ©Carlos Guestrin 2006

30

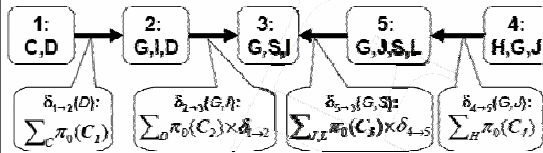
Choice of root

- Message does not depend on root!!!

Root: node 5

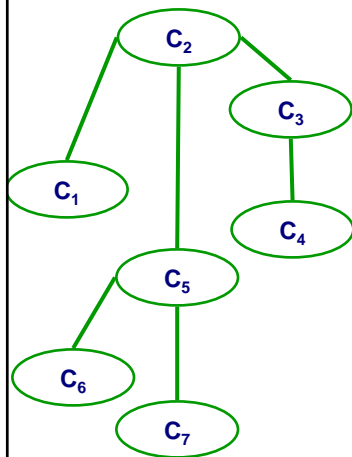


Root: node 3



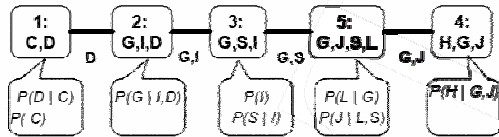
“Cache” computation: Obtain belief for all roots in linear time!!

Shafer-Shenoy Algorithm (a.k.a. VE in clique tree for all roots)



- Clique C_i ready to transmit to neighbor C_j if received messages from all neighbors but j
 - Leaves are always ready to transmit
- While $\exists C_i$ ready to transmit to C_j
 - Send message $\delta_{i \rightarrow j}$
- Complexity: Linear in # cliques
 - One message sent each direction in each edge
- Corollary: At convergence
 - Every clique has correct belief

Calibrated Clique tree



- Initially, neighboring nodes don't agree on "distribution" over separators
- **Calibrated clique tree:**
 - At convergence, tree is *calibrated*
 - Neighboring nodes agree on distribution over separator

10.708 – ©Carlos Guestrin 2006

33

Answering queries with clique trees

- Query within clique
- Incremental updates – Observing evidence $Z=z$
 - Multiply some clique by indicator $\mathbf{1}(Z=z)$
- Query outside clique
 - Use variable elimination!

10.708 – ©Carlos Guestrin 2006

34

Message passing with division



- Computing messages by multiplication:

- Computing messages by division:

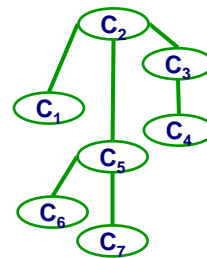
10.708 – ©Carlos Guestrin 2006

35

Lauritzen-Spiegelhalter Algorithm (a.k.a. belief propagation)

Simplified description
see reading for details

- Initialize all separator potentials to 1
 - $\mu_{ij} \leftarrow 1$
- All messages ready to transmit
- While $\exists \delta_{i \rightarrow j}$ ready to transmit
 - $\mu_{ij}' \leftarrow$
 - If $\mu_{ij}' \neq \mu_{ij}$
 - $\delta_{i \rightarrow j} \leftarrow$
 - $\pi_j \leftarrow \pi_j \times \delta_{i \rightarrow j}$
 - $\mu_{ij} \leftarrow \mu_{ij}'$
 - \forall neighbors k of j , $k \neq i$, $\delta_{j \rightarrow k}$ ready to transmit
- Complexity: Linear in # cliques
 - for the “right” schedule over edges (leaves to root, then root to leaves)
- **Corollary:** At convergence, every clique has correct belief



10.708 – ©Carlos Guestrin 2006

36

VE versus BP in clique trees

- VE messages (the one that multiplies)

- BP messages (the one that divides)

10-708 - ©Carlos Guestrin 2006

37

Clique tree invariant

- **Clique tree potential:**
 - Product of clique potentials divided by separators potentials

- **Clique tree invariant:**
 - $P(\mathbf{X}) = \pi_T(\mathbf{X})$

10-708 - ©Carlos Guestrin 2006

38

Belief propagation and clique tree invariant

- **Theorem:** Invariant is maintained by BP algorithm!

- BP reparameterizes clique potentials and separator potentials
 - At convergence, potentials and messages are marginal distributions

10-708 – ©Carlos Guestrin 2006

39

Subtree correctness

- **Informed message** from i to j , if all messages into i (other than from j) are informed
 - Recursive definition (leaves always send informed messages)
- **Informed subtree:**
 - All incoming messages informed
- **Theorem:**
 - Potential of connected informed subtree T' is marginal over $\text{scope}[T']$
- **Corollary:**
 - At convergence, clique tree is *calibrated*
 - $\pi_i = P(\text{scope}[\pi_i])$
 - $\mu_{ij} = P(\text{scope}[\mu_{ij}])$

10-708 – ©Carlos Guestrin 2006

40

Clique trees versus VE

- Clique tree advantages
 - Multi-query settings
 - Incremental updates
 - Pre-computation makes complexity explicit

- Clique tree disadvantages
 - Space requirements – no factors are “deleted”
 - Slower for single query
 - Local structure in factors may be lost when they are multiplied together into initial clique potential

10-708 – ©Carlos Guestrin 2006

41

Clique tree summary

- Solve marginal queries for all variables in only twice the cost of query for one variable
- Cliques correspond to maximal cliques in induced graph
- Two message passing approaches
 - VE (the one that multiplies messages)
 - BP (the one that divides by old message)
- Clique tree invariant
 - Clique tree potential is always the same
 - We are only reparameterizing clique potentials
- Constructing clique tree for a BN
 - from elimination order
 - from triangulated (chordal) graph
- Running time (only) exponential in size of largest clique
 - Solve **exactly** problems with thousands (or millions, or more) of variables, and cliques with tens of nodes (or less)

10-708 – ©Carlos Guestrin 2006

42