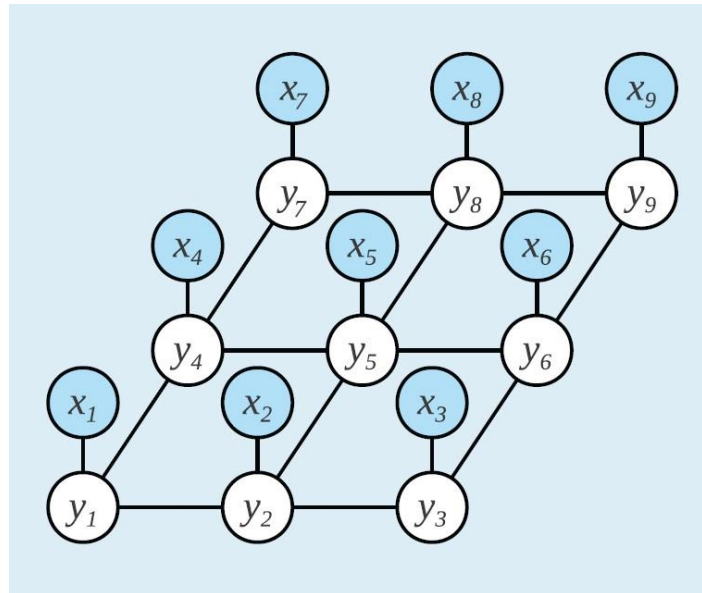


Probabilistic Graphical Models in Bioinformatics

Lecture 13: Clique trees; MAP inference



What you have learned so far **and remaining topics**

- **Representation**

- Bayesian networks
- Markov networks
- Conditional independencies; conditional probability distributions; continuous and discrete data
- **Temporal models: HMM and dynamic Bayesian networks**

- **Inference**

- Variable elimination
- **Clique trees**
- **MAP inference**
- **Sampling methods: Gibbs sampling and MCMC**

- **Learning**

- Parameter learning: MLE, Bayesian inference
- Structure learning: constraint-based and score-based methods
- In the case of partially observed data: gradient ascent methods; EM

- **Applications**

- Gene regulatory networks
- Motif finding
- **Profile HMM for sequence alignment**

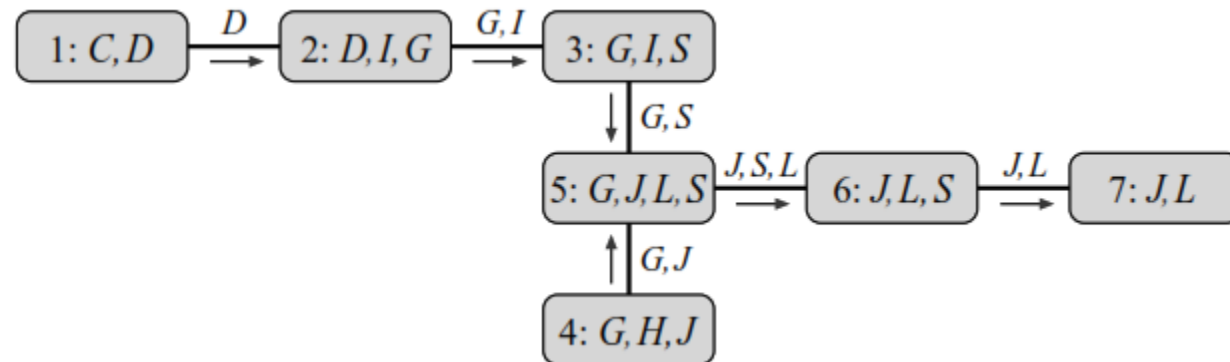
Exact inference: clique tree

Exact inference: clique tree

- We already showed how variable elimination can perform exact inference effectively.
 - Key insight: performing local operations on the factors defining the distribution, rather than generating the entire joint distribution.
- Today we present an alternative implementation of the same insight.
 - The algorithm uses a global data structure, called clique tree, for scheduling manipulations of factors.
- In the variable elimination,
 - each step in the computation creates a factor ψ_i by multiplying existing factors.
 - A variable is then eliminated in ψ_i to generate a new factor τ_i .
- In this section, we consider a factor ψ_i to be a computational data structure
 - takes “messages” τ_j generated by other factor ψ_j
 - and generates a message τ_i used by another factor ψ_l .

Cluster tree definition

- A cluster graph U for a set of factors Φ over \mathcal{X} is an undirected graph
 - each of whose nodes i is associated with a subset $C_i \subset \mathcal{X}$
 - each edge between a pair of clusters C_i and C_j is associated with a sepset $S_{i,j} \subseteq C_i \cap C_j$.
- A cluster graph must be family-preserving
 - Each factor $\phi \in \Phi$ must be associated with a cluster C_i

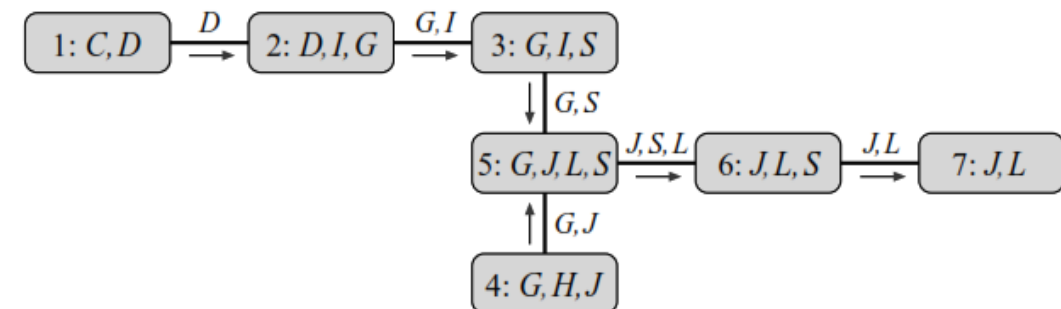


Variable elimination and cluster graph

Example revisited

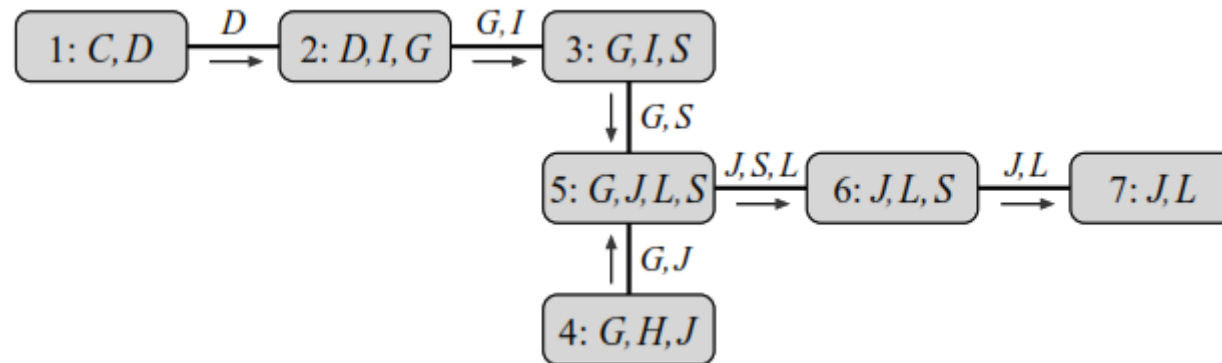
- We have seven factors ψ_1, \dots, ψ_7 .
- The message $\tau_1(D)$ is generated from $\psi_1(C, D)$.
- $\tau_1(D)$ participates in the computation of ψ_2 .
- Hence, there is an edge between C_1 and C_2 .

Step	Variable eliminated	Factors used	Variables involved	New factor
1	C	$\phi_C(C), \phi_D(D, C)$	C, D	$\tau_1(D)$
2	D	$\phi_G(G, I, D), \tau_1(D)$	G, I, D	$\tau_2(G, I)$
3	I	$\phi_I(I), \phi_S(S, I), \tau_2(G, I)$	G, S, I	$\tau_3(G, S)$
4	H	$\phi_H(H, G, J)$	H, G, J	$\tau_4(G, J)$
5	G	$\tau_4(G, J), \tau_3(G, S), \phi_L(L, G)$	G, J, L, S	$\tau_5(J, L, S)$
6	S	$\tau_5(J, L, S), \phi_J(J, L, S)$	J, L, S	$\tau_6(J, L)$
7	L	$\tau_6(J, L)$	J, L	$\tau_7(J)$



Running intersection property

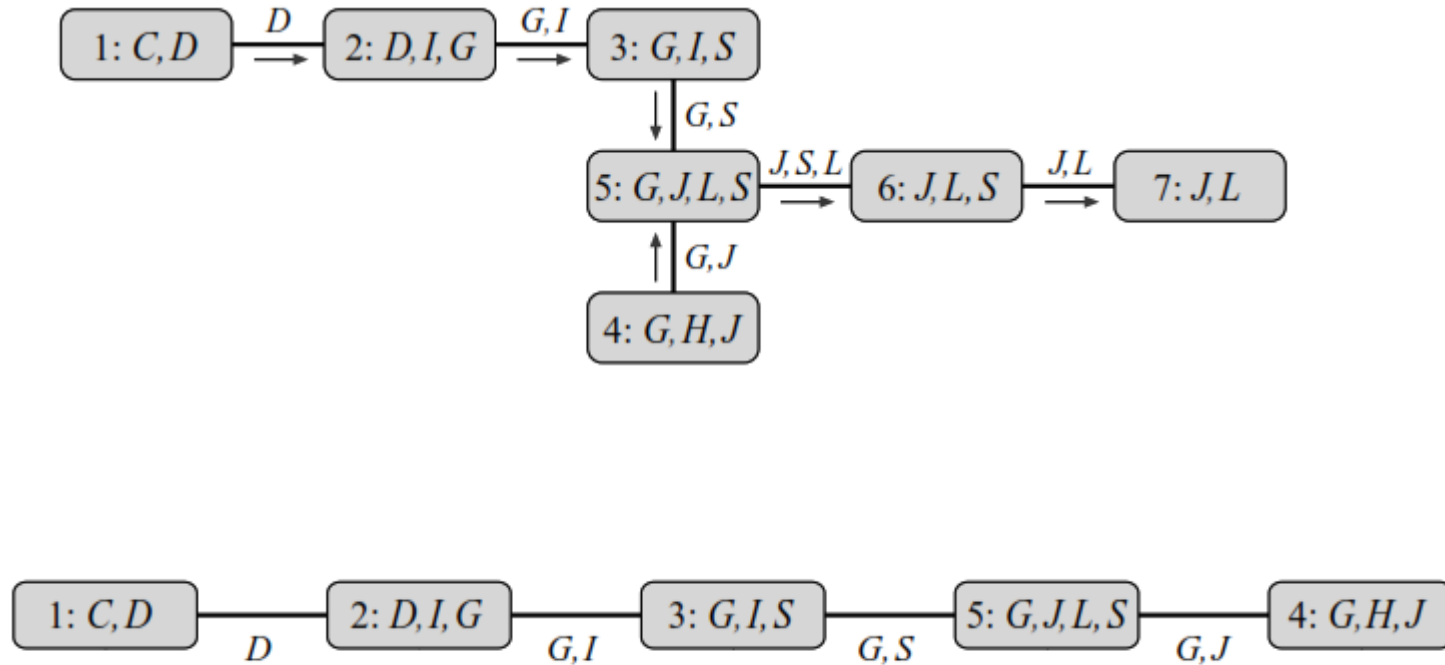
- We say a cluster tree has the running intersection property if $X \in C_i$ and $X \in C_j$, then X is also in **every** cluster in the **unique** path between C_i and C_j .



Note that running intersection property implies that $S_{i,j} = C_i \cap C_j$.

Clique tree

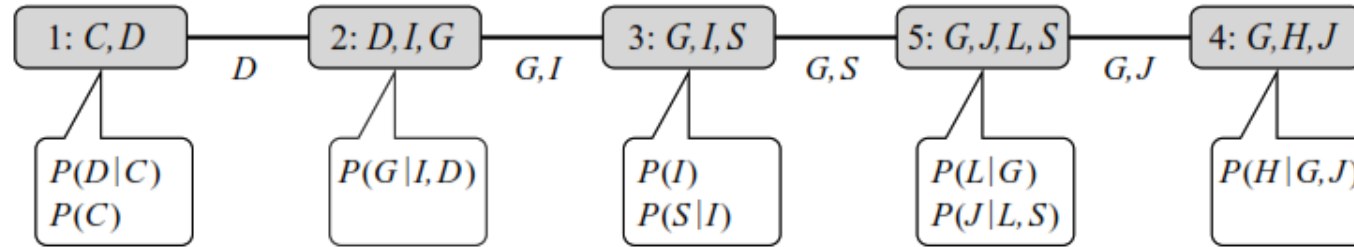
- A cluster tree that satisfies the running intersection property is called a clique tree, also called a junction tree or a joint tree.



Nonmaximal cliques are absent in this simplified clique tree.

Variable elimination in a clique tree

- Consider one possible clique tree for the Student network



- The first step is to generate a set of *initial potentials* associated with the different cliques by multiplying the initial factors assigned each clique

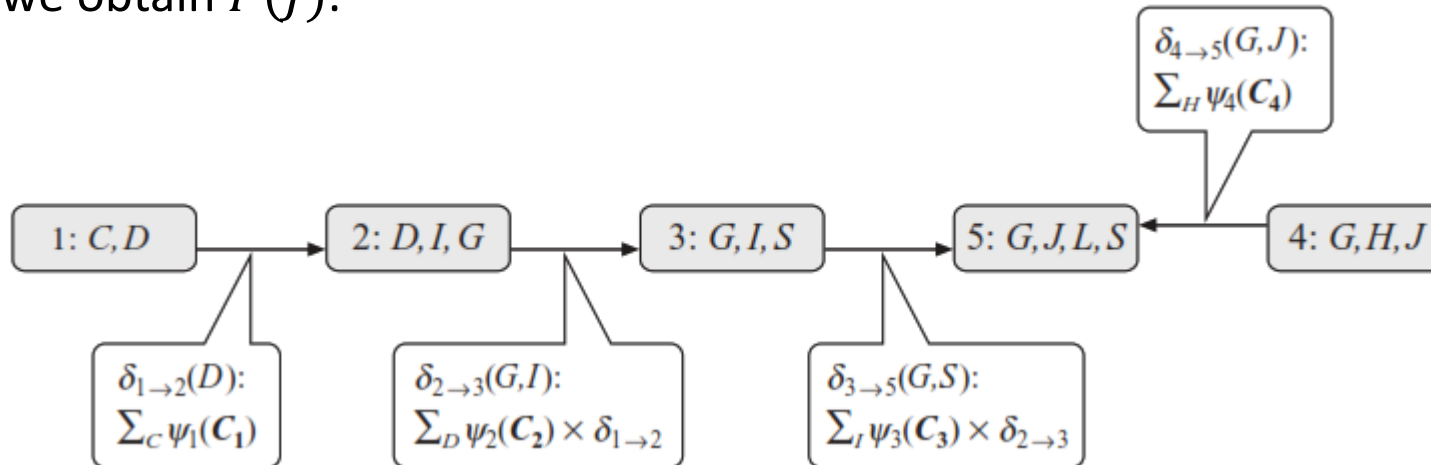
$$\psi_5(J, L, G, S) = \phi_L(L, G) \cdot \phi_J(J, L, S)$$

- Assume our task is to compute $P(J)$. Then we want to do VE so the J is not eliminated.
- We select as our root clique some clique that contains J , for example C5

Variable elimination in a clique tree

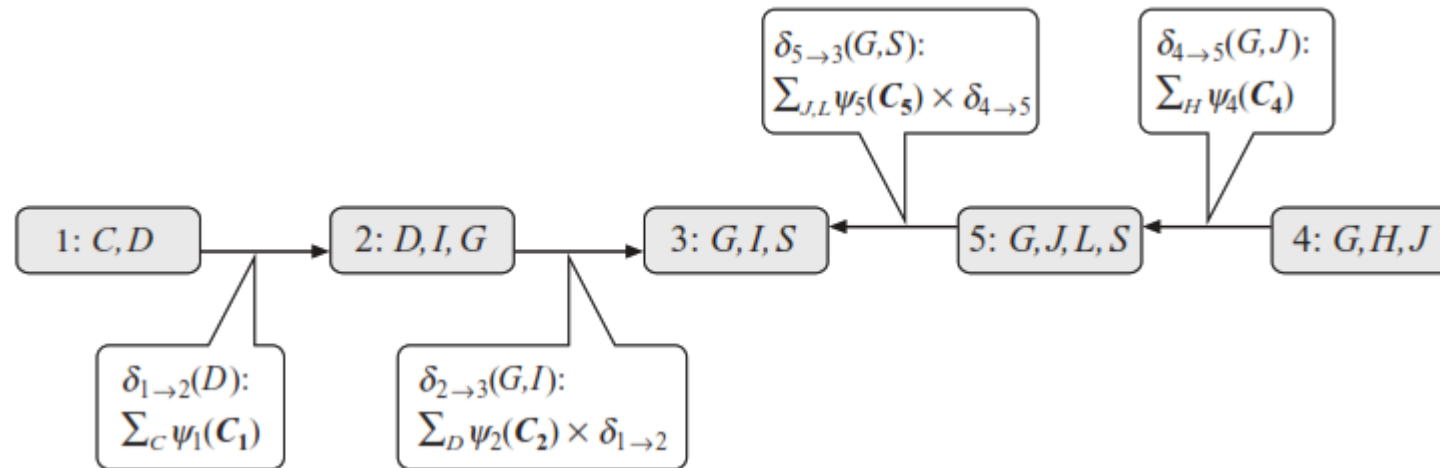
- We then execute the following steps:

1. In C_1 : we eliminate C by performing $\sum_C \psi_1(C, D)$. The resulting factor has scope D and send it as a message $\delta_{1 \rightarrow 2}(D)$ to C_2 .
2. In C_2 : we eliminate D . The resulting factor has scope G, I and send it as a message $\delta_{2 \rightarrow 3}(G, I)$ to C_3 .
3. In C_3 : we eliminate I . The resulting factor has scope G, S and send it as a message $\delta_{3 \rightarrow 5}(G, S)$ to C_5 .
4. In C_4 : we eliminate H . The resulting factor has scope G, J and send it as a message $\delta_{4 \rightarrow 5}(G, J)$ to C_5 .
5. In C_5 : we define $\beta_5(G, J, S, L) = \delta_{3 \rightarrow 5}(G, S) \cdot \delta_{4 \rightarrow 5}(G, J) \cdot \psi_5(G, J, S, L)$. By summing out G, L , and S we obtain $P(J)$.



Another example: computing $P(G)$

- Messages sent from C_1 to C_2 , from C_2 to C_3 , and from C_4 to C_5 are the same as the previous run



- Ready clique: a clique is ready when it has received all of its incoming messages.

Clique-tree message passing

1. Let T be a clique tree with the cliques C_1, \dots, C_k .
2. We begin by multiplying the factors assigned to each clique, resulting in our initial potentials.

Because each factor is assigned to exactly one clique, we have

$$\psi_j(C_j) = \prod_{\phi : \alpha(\phi)=j} \phi. \quad \longrightarrow \quad \prod_{\phi} \phi = \prod_j \psi_j.$$

3. We use clique-tree data structure to pass messages between neighboring cliques, sending all messages toward to the root clique.

$$\delta_{i \rightarrow j} = \sum_{C_i - S_{i,j}} \psi_i \cdot \prod_{k \in (\text{Nb}_i - \{j\})} \delta_{k \rightarrow i}.$$

4. When the root clique has received all messages, it multiplies them with its own initial potential. The result is a factor called the *beliefs*, denoted $\beta_r(C_r)$. It represents

$$\tilde{P}_{\Phi}(C_r) = \sum_{\mathcal{X} - C_r} \prod_{\phi} \phi.$$

Algorithm 10.1 Upward pass of variable elimination in clique tree

```

Procedure CTree-SP-Upward (
     $\Phi$ ,    // Set of factors
     $\mathcal{T}$ ,    // Clique tree over  $\Phi$ 
     $\alpha$ ,    // Initial assignment of factors to cliques
     $C_r$     // Some selected root clique
)
1  Initialize-Cliques
2  while  $C_r$  is not ready
3      Let  $C_i$  be a ready clique
4       $\delta_{i \rightarrow p_r(i)}(S_{i, p_r(i)}) \leftarrow \text{SP-Message}(i, p_r(i))$ 
5       $\beta_r \leftarrow \psi_r \cdot \prod_{k \in \text{Nb}_{C_r}} \delta_{k \rightarrow r}$ 
6  return  $\beta_r$ 

```

```

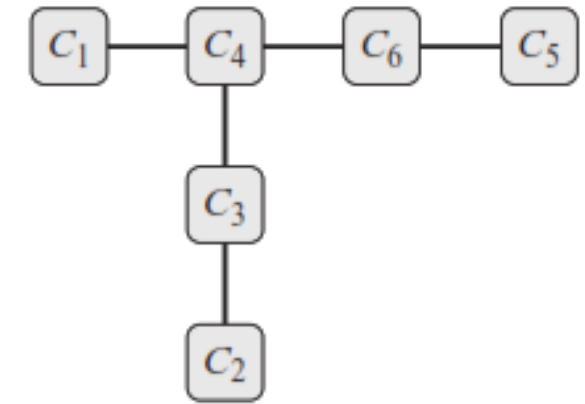
Procedure Initialize-Cliques (
)
1  for each clique  $C_i$ 
2       $\psi_i(C_i) \leftarrow \prod_{\phi_j : \alpha(\phi_j)=i} \phi_j$ 
3

```

```

Procedure SP-Message (
    i,    // sending clique
    j     // receiving clique
)
1   $\psi(C_i) \leftarrow \psi_i \cdot \prod_{k \in (\text{Nb}_i - \{j\})} \delta_{k \rightarrow i}$ 
2   $\tau(S_{i,j}) \leftarrow \sum_{C_i - S_{i,j}} \psi(C_i)$ 
3  return  $\tau(S_{i,j})$ 

```

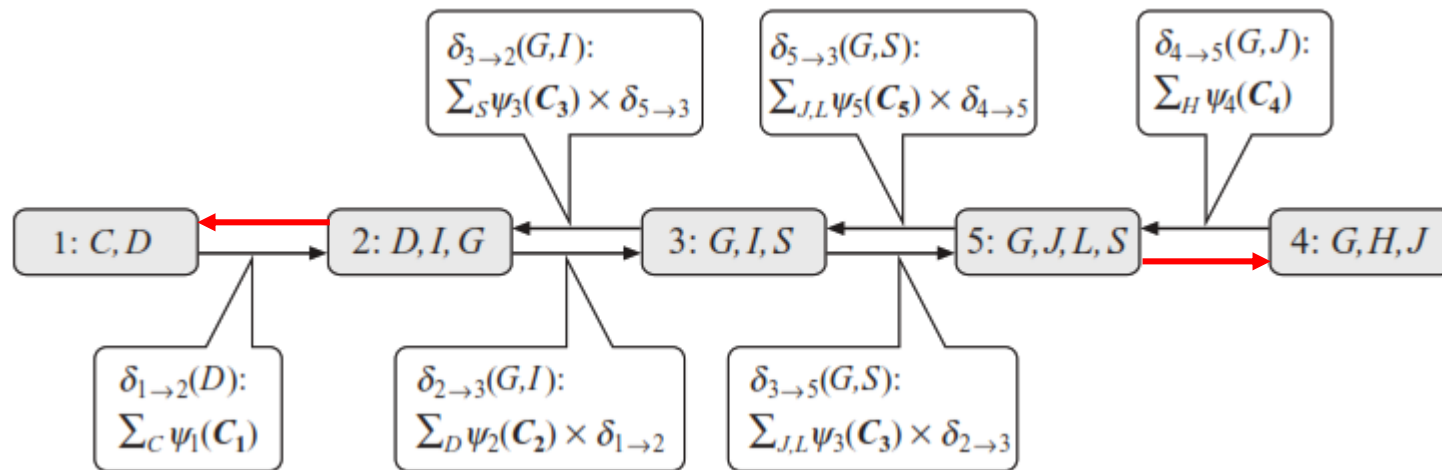


*Assume we have selected C_6 as our root clique.
There are multiple legitimate orderings.*

Clique tree calibration

- We showed how to use the same clique tree to compute the probability of any variable in the input variables.
- We often wish to estimate the probability of a large number of variables.
 - For example, we want to compute the probability of several possible diseases in medical diagnosis setting.
- Consider the task of computing the posterior distribution over each random variable in the network
 - Naïve approach: doing inference separately for each variable.
 - Better approach: ***sum-product belief propagation*** algorithm.

Sum-Product Belief Propagation



Sum-Product Belief Propagation

Algorithm 10.2 Calibration using sum-product message passing in a clique tree

```
Procedure CTree-SP-Calibrate (  
     $\Phi$ , // Set of factors  
     $\mathcal{T}$  // Clique tree over  $\Phi$   
)  
1  Initialize-Cliques  
2  while exist  $i, j$  such that  $i$  is ready to transmit to  $j$   
3       $\delta_{i \rightarrow j}(S_{i,j}) \leftarrow \text{SP-Message}(i, j)$   
4  for each clique  $i$   
5       $\beta_i \leftarrow \psi_i \cdot \prod_{k \in \text{Nb}_i} \delta_{k \rightarrow i}$   
6  return  $\{\beta_i\}$ 
```

Scheduling for sending messages:

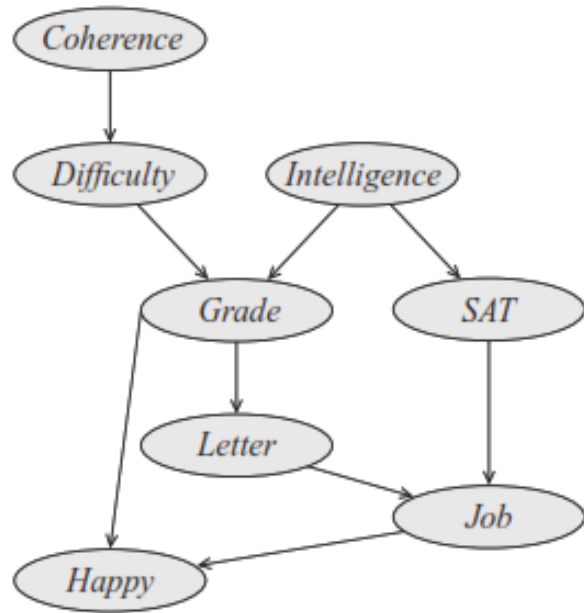
- Asynchronously
- Upward pass and downward pass

Constructing a clique tree

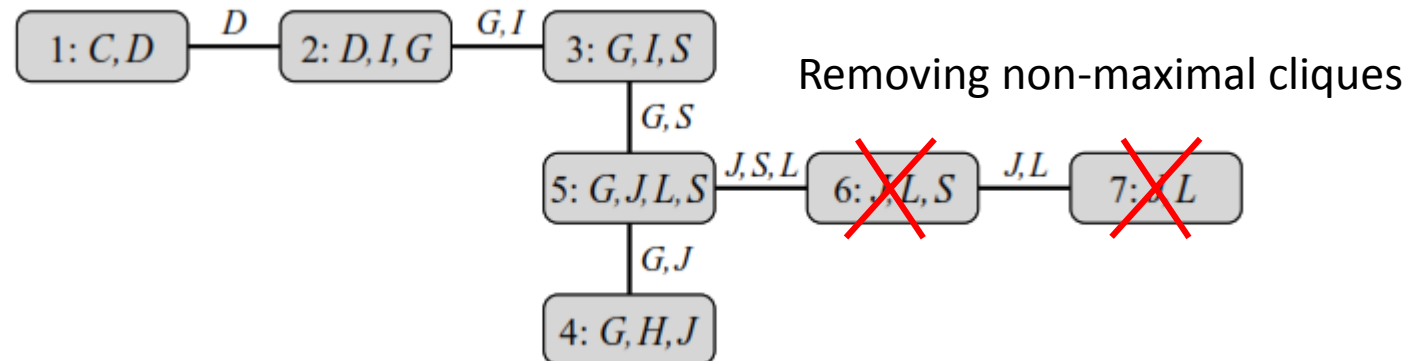
Constructing a clique tree

- How do we construct a clique tree for a set of factors, or equivalently, for its underlying undirected graph H_Φ ? (*See section 9.4.2.1 for the definition of H_Φ .*)
- Two basic approaches
 - Clique trees from variable elimination
 - Clique trees from chordal graphs

Clique trees from variable elimination



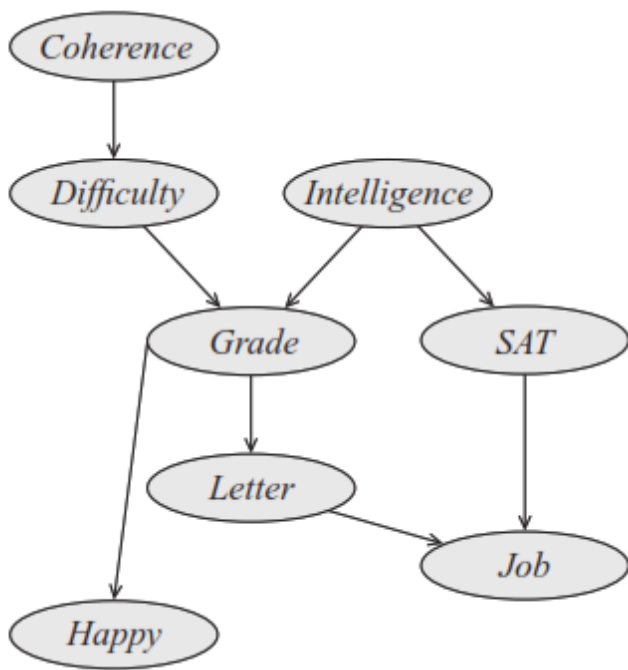
Step	Variable eliminated	Factors used	Variables involved	New factor
1	C	$\phi_C(C), \phi_D(D, C)$	C, D	$\tau_1(D)$
2	D	$\phi_G(G, I, D), \tau_1(D)$	G, I, D	$\tau_2(G, I)$
3	I	$\phi_I(I), \phi_S(S, I), \tau_2(G, I)$	G, S, I	$\tau_3(G, S)$
4	H	$\phi_H(H, G, J)$	H, G, J	$\tau_4(G, J)$
5	G	$\tau_4(G, J), \tau_3(G, S), \phi_L(L, G)$	G, J, L, S	$\tau_5(J, L, S)$
6	S	$\tau_5(J, L, S), \phi_J(J, L, S)$	J, L, S	$\tau_6(J, L)$
7	L	$\tau_6(J, L)$	J, L	$\tau_7(J)$



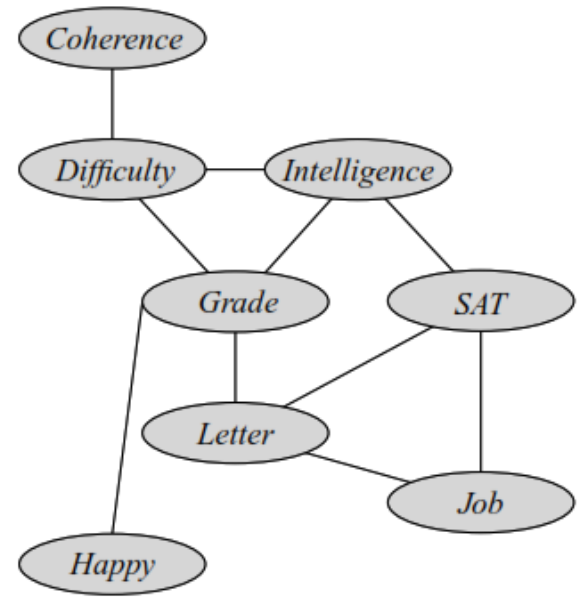
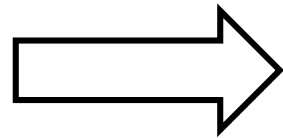
- Each cluster is defined by $C_i = \text{Scope}[\psi_i]$
- There is an edge between C_i and C_j if the message τ_i is used in computing ψ_j

Clique trees from chordal graphs

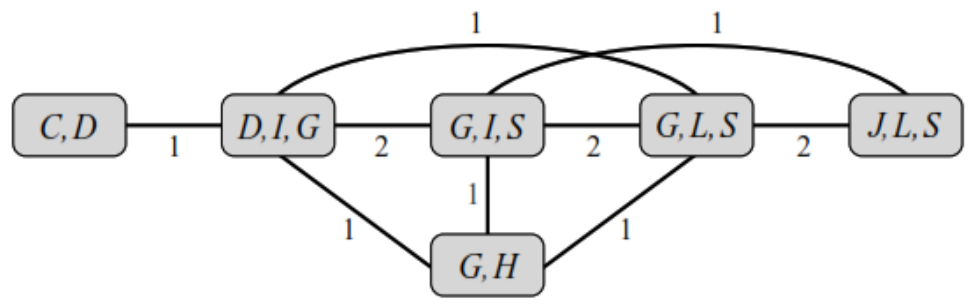
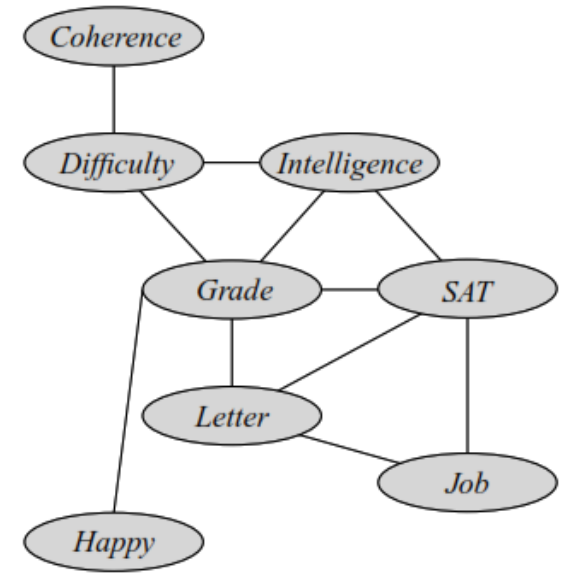
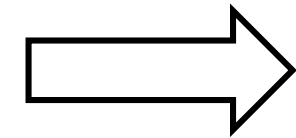
1. Find the minimal Markov network I-map H_Φ
 - For a Bayesian network, in the case without evidence, it is simply the moralized graph
2. Triangulation to obtain a chordal graph H^*
3. Find cliques in H^* , and make each one a node in the cluster graph
4. Run the maximum spanning tree algorithm on the cluster graph to construct a tree



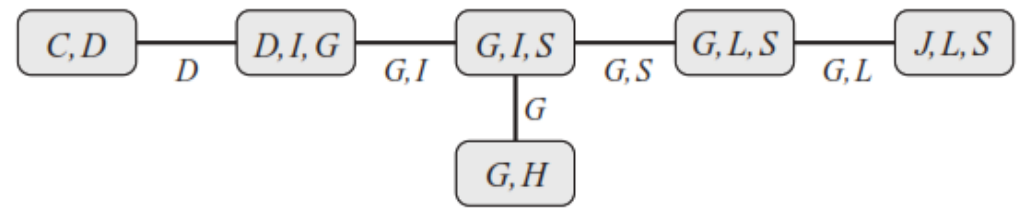
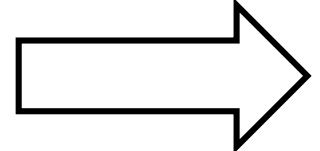
moralization



triangulation



Maximum spanning tree



MAP inference

MAP inference

- A MAP query aims to find the most likely assignment to all of the (non-evidence) variables.
- A marginal MAP query aims to find the most likely assignment to a subset of the variables, marginalizing over the rest.
- Examples:
 - Speech recognition: to decode the most likely utterance given the (noisy) acoustic signal.
- For addressing MAP queries, we need to compute

$$\xi^{map} = \arg \max_{\xi} P_{\Phi}(\xi) = \arg \max_{\xi} \frac{1}{Z} \tilde{P}_{\Phi}(\xi) = \arg \max_{\xi} \tilde{P}_{\Phi}(\xi).$$

Max-Product Variable Elimination

- Consider the Bayesian network $A \rightarrow B$. Assume our goal is to compute:

$$\begin{aligned}\max_{a,b} P(a,b) &= \max_{a,b} P(a)P(b | a) \\ &= \max_a \max_b P(a)P(b | a).\end{aligned}$$

- For any choice of a , the value of B must be chosen so as to maximize $P(b | a)$. We must also choose the value of A appropriately

$$\max_b P(a)P(b | a) = P(a) \max_b P(b | a).$$

- Let $\phi(a)$ denote $\max_b P(b | a)$. For example, for the following assignments

a^0	a^1	A	b^0	b^1
0.4	0.6	a^0	0.1	0.9
		a^1	0.55	0.45

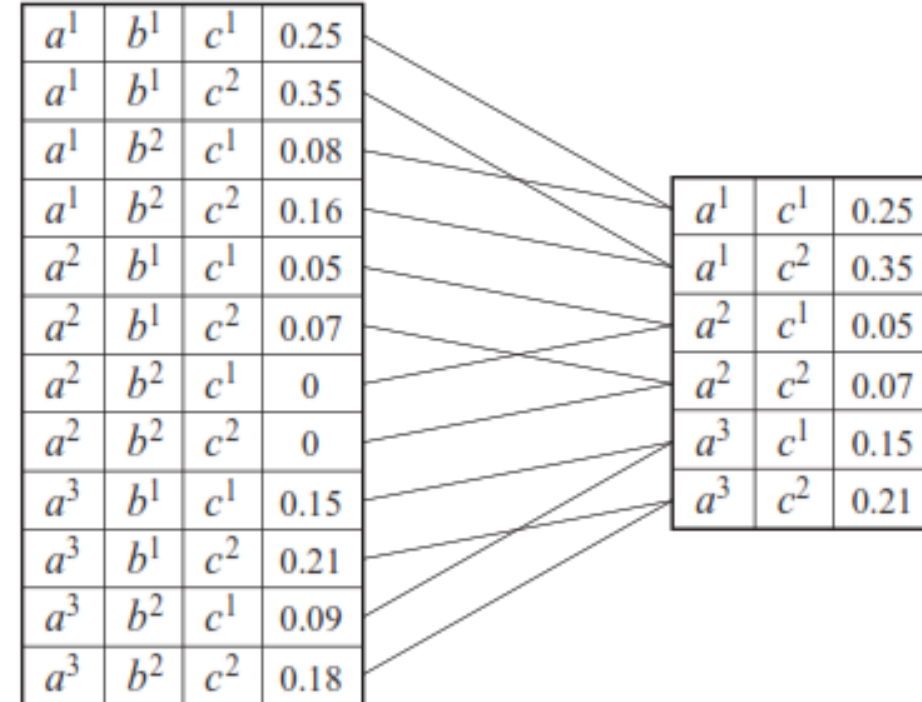
- $\phi(a^1) = \max_b P(b | a^1) = 0.55$ and $\phi(a^0) = \max_b P(b | a^0) = 0.9$. Finally

$$\max_a P(a)\phi(a) = \max [0.4 \cdot 0.9, 0.6 \cdot 0.55] = 0.36.$$

Max-marginalization

Let $\phi(X, Y)$ be a factor. We define the factor maximization of Y in ϕ as

$$\psi(\mathbf{X}) = \max_Y \phi(\mathbf{X}, Y).$$



Algorithm 13.1 Variable elimination algorithm for MAP. The algorithm can be used both in its max-product form, as shown, or in its max-sum form, replacing factor product with factor addition.

```

Procedure Max-Product-VE (
     $\Phi$ , // Set of factors over  $\mathbf{X}$ 
     $\prec$  // Ordering on  $\mathbf{X}$ 
)
1  Let  $X_1, \dots, X_k$  be an ordering of  $\mathbf{X}$  such that
2   $X_i \prec X_j$  iff  $i < j$ 
3  for  $i = 1, \dots, k$ 
4   $(\Phi, \phi_{X_i}) \leftarrow \text{Max-Product-Eliminate-Var}(\Phi, X_i)$ 
5   $\mathbf{x}^* \leftarrow \text{Traceback-MAP}(\{\phi_{X_i} : i = 1, \dots, k\})$ 
6  return  $\mathbf{x}^*, \Phi$  //  $\Phi$  contains the probability of the MAP

```

```

Procedure Max-Product-Eliminate-Var (
     $\Phi$ , // Set of factors
     $Z$  // Variable to be eliminated
)
1   $\Phi' \leftarrow \{\phi \in \Phi : Z \in \text{Scope}[\phi]\}$ 
2   $\Phi'' \leftarrow \Phi - \Phi'$ 
3   $\psi \leftarrow \prod_{\phi \in \Phi'} \phi$ 
4   $\tau \leftarrow \max_Z \psi$ 
5  return  $(\Phi'' \cup \{\tau\}, \psi)$ 

```

```

Procedure Traceback-MAP (
     $\{\phi_{X_i} : i = 1, \dots, k\}$ 
)
1  for  $i = k, \dots, 1$ 
2   $\mathbf{u}_i \leftarrow (x_{i+1}^*, \dots, x_k^*) \langle \text{Scope}[\phi_{X_i}] - \{X_i\} \rangle$ 
3  // The maximizing assignment to the variables eliminated after
    $X_i$ 
4   $x_i^* \leftarrow \arg \max_{x_i} \phi_{X_i}(x_i, \mathbf{u}_i)$ 
5  //  $x_i^*$  is chosen so as to maximize the corresponding entry in
   the factor, relative to the previous choices  $\mathbf{u}_i$ 
6  return  $\mathbf{x}^*$ 

```

Theorem 13.4

The algorithm of algorithm 13.1 returns

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} \prod_{\phi \in \Phi} \phi,$$

and Φ , which contains a single factor of empty scope whose value is:

$$\max_{\mathbf{x}} \prod_{\phi \in \Phi} \phi.$$

A run of max-product variable elimination



Step	Variable eliminated	Factors used	Intermediate factor	New factor
1	S	$\phi_S(I, S)$	$\psi_1(I, S)$	$\tau_1(I)$
2	I	$\phi_I(I), \phi_G(G, I, D), \tau_1(I)$	$\psi_2(G, I, D)$	$\tau_2(G, D)$
3	D	$\phi_D(D), \tau_2(G, D)$	$\psi_3(G, D)$	$\tau_3(G)$
4	L	$\phi_L(L, G)$	$\psi_4(L, G)$	$\tau_4(G)$
5	G	$\tau_4(G), \tau_3(G)$	$\psi_5(G)$	$\tau_5(\emptyset)$

For example, the first step would compute

$$\tau_1(I) = \max_s \phi_S(I, s)$$

The final factor, $\tau_5(\emptyset)$, is simply a number, whose value is

$$\max_{S, I, D, L, G} P(S, I, D, L, G).$$

Decoding or finding the most probable assignment

- Traceback of the solution
- We begin by computing g

$$g^* = \arg \max_g \psi_5(g)$$

$$l^* = \arg \max_l \psi_4(g^*, l)$$

$$d^* = \arg \max_d \psi_3(g^*, d)$$

$$i^* = \arg \max_i \psi_2(g^*, i, d^*)$$

$$s^* = \arg \max_s \psi_1(i^*, s).$$

Step	Variable eliminated	Factors used	Intermediate factor	New factor
1	S	$\phi_S(I, S)$	$\psi_1(I, S)$	$\tau_1(I)$
2	I	$\phi_I(I), \phi_G(G, I, D), \tau_1(I)$	$\psi_2(G, I, D)$	$\tau_2(G, D)$
3	D	$\phi_D(D), \tau_2(G, D)$	$\psi_3(G, D)$	$\tau_3(G)$
4	L	$\phi_L(L, G)$	$\psi_4(L, G)$	$\tau_4(G)$
5	G	$\tau_4(G), \tau_3(G)$	$\psi_5(G)$	$\tau_5(\emptyset)$

(Refer to Example 13.4 for more details)