

Constructing a clique tree from VE



- Select elimination order <</p>
- Connect factors that would be generated if you run VE with order ≺
- Simplify!
 - □ Eliminate factor that is subset G7∠ G7 of neighbor

10-708 - ©Carlos Guestrin 2006

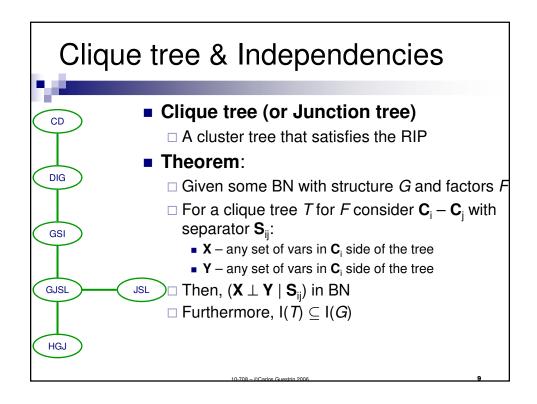
Find clique tree from chordal graph

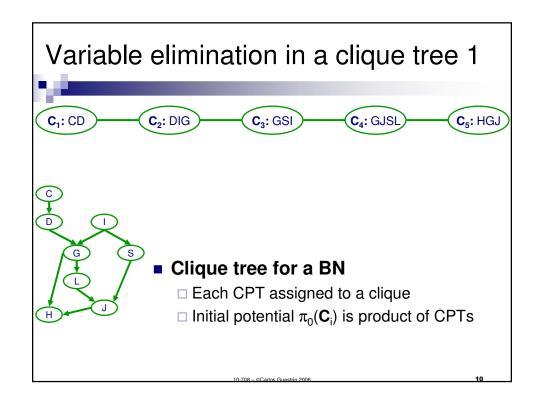


- Triangulate moralized graph to obtain chordal graph
- Find maximal cliques
 - □ NP-complete in general
 - $\ \square$ Easy for chordal graphs
 - □ Max-cardinality search
- Maximum spanning tree finds clique tree satisfying RIP!!!
 - ☐ Generate weighted graph over cliques
 - □ Edge weights (i,j) is separator size |C_i∩C_i|

Coherence
Difficulty
Intelligence
Grade
SAT
Letter
Job

10-708 - ©Carlos Guestrin 2006





Variable elimination in a clique tree 2



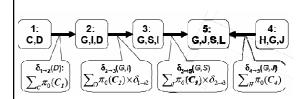
■ VE in clique tree to compute P(X_i)

- □ 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

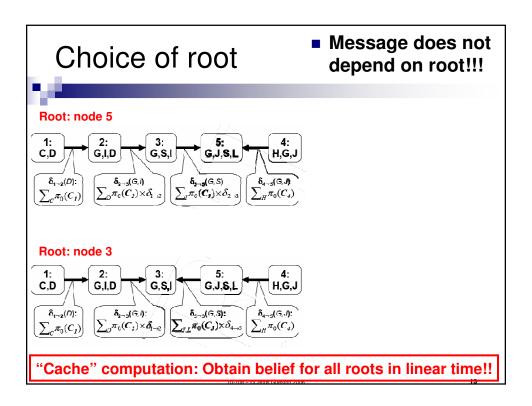
44

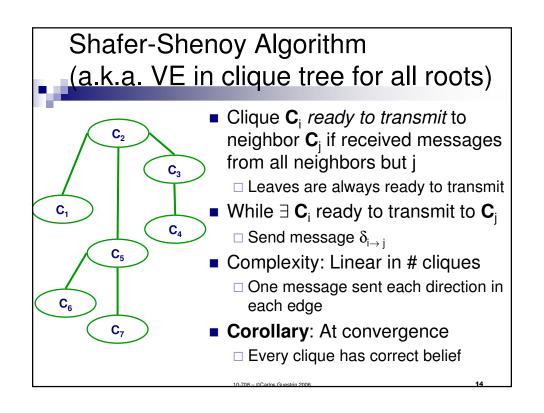
Belief from message



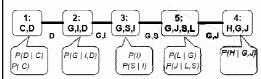
- Theorem: When clique C_i is ready
 - □ Received messages from all neighbors
 - \square Belief $\pi_i(\mathbf{C}_i)$ is product of initial factor with messages:

10-708 = @Carlos Guestrin 2006





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

15

Answering queries with clique trees



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

10-708 - ©Carlos Guestrin 2006

Message passing with division



- Computing messages by multiplication:
- Computing messages by division:

10-708 = @Carlos Guestrin 2006

17

Lauritzen-Spiegelhalter Algorithm

(a.k.a. belief propagation)

Simplified description see reading for details

- Initialize all separator potentials to 1
 - $\square \ \mu_{ij} \leftarrow 1$
- All messages ready to transmit
- \blacksquare While $\exists~\delta_{i\rightarrow~j}$ ready to transmit
 - $\square \; \mu_{ii} \; \dot{} \leftarrow$
 - \square If μ_{ii} $\neq \mu_{ii}$
 - $\delta_{i \rightarrow i} \leftarrow$
 - $\quad \blacksquare \quad \pi_j \leftarrow \pi_j \times \delta_{i \rightarrow j}$
 - $\quad \blacksquare \quad \mu_{ii} \leftarrow \mu_{ii},$
 - $\quad \forall$ neighbors k of j, k≠ 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

VE versus BP in clique trees



- VE messages (the one that multiplies)
- BP messages (the one that divides)

10-708 - ©Carlos Guestrin 2006

19

Clique tree invariant



- Clique tree potential:
 - □ Product of clique potentials divided by separators potentials
- Clique tree invariant:
 - $\square P(\mathbf{X}) = \pi_T(\mathbf{X})$

I0-708 – ©Carlos Guestrin 2006

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

**O-rice Overhile 0000

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:
 - $\hfill \square$ All incoming messages informed
- Theorem:
 - □ Potential of connected informed subtree *T'* is marginal over scope[*T'*]
- Corollary:
 - ☐ At convergence, clique tree is *calibrated*
 - $\pi_i = P(scope[\pi_i])$
 - $\mu_{ij} = P(scope[\mu_{ij}])$

8 = @Carlos Guestrin 2006

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

23

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

Announcements



- Recitation tomorrow, don't miss it!!!
 - □ Ajit on Junction Trees

10-708 - ©Carlos Guestrin 2006

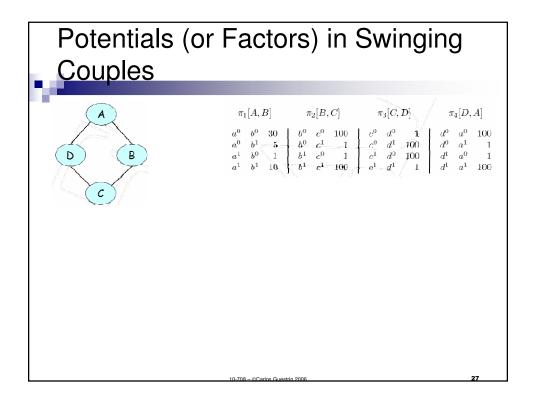
25

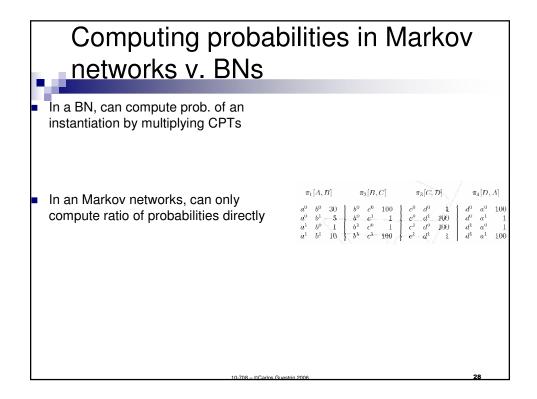
Swinging Couples revisited



- This is no perfect map in BNs
- But, an undirected model will be a perfect map

)-708 = ©Carlos Guestrin 2006



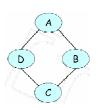


Normalization for computing probabilities

 To compute actual probabilities, must compute normalization constant (also called partition function)

1	Assignment				Unnormalized	Noncolinad
ı					L 4 \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \	
ı	a^0	b^0	c^{o}	d^0	300000	0.071
ı	a^0	b^0	e^0	d^1	300000	0.04
ı	a^0	b^{0}	c^{\perp}	d^0	300000	0.04
ı	a^0	b^{Ω}	e^{1}	d^{\dagger}	30	$4.1 \cdot 10^{-6}$
ı	a^0	b^3	e^{0}	136	500	6.9 - 10 - 5
١	a^0	b^i	$c^{\scriptscriptstyle \mathrm{D}}$	d^1	590	6.9 - 10 - 5
١	a^0	b^1	c^1	d^0	5000000	0.69
١	a^0	b^1	c^{1}	d^{1}	500	6.9 · 10 5
١	a^1	b^0	v^{0}	420	100	1.4 - 18) 5
١	a^{1}	b^{o}	ęa.	d^{\perp}	1000000	≥0.14
ı	a^1	b^{Ω}	e^{i}	d^{α}	100	1.4 .10 5
ı	a^{1}	b^{0}	c^{1}	d^{1}	190	$1.4 \cdot 10^{-5}$
ı	a^1	b1	e^{α}	d^{0}	10	1.4 - 10 6
ı	a^1	b^1	e^0	d^1	100000	0.014
ı	a^1	b^1	c^1	d^0	1000000	0.014
ì	a^{i}	b^{\pm}	c^{i}	$d^{\underline{\epsilon}}$	100000	0.014

 \blacksquare Computing partition function is hard! \to Must sum over all possible assignments



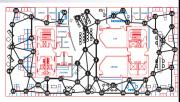
10-708 - ©Carlos Guestrin 2006

20

Factorization in Markov networks



- Given an undirected graph *H* over variables **X**={X₁,...,X_n}
- A distribution *P* **factorizes** over *H* if ∃
 - $\ \square$ subsets of variables $\mathbf{D_1} \subseteq \mathbf{X},...,\,\mathbf{D_m} \subseteq \mathbf{X},$ such that the $\mathbf{D_i}$ are *fully connected* in H
 - $\hfill\Box$ non-negative potentials (or factors) $\pi_{_{1}}(\boldsymbol{D_{1}}), \ldots, \, \pi_{_{m}}(\boldsymbol{D_{m}})$
 - also known as clique potentials
 - such that



 Also called Markov random field H, or Gibbs distribution over H

10-708 - ©Carlos Guestrin 2006

Global Markov assumption in Markov networks



A path $X_1 - ... - X_k$ is **active** when set of variables **Z** are observed if none of $X_i \in \{X_1,...,X_k\}$ are observed (are part of **Z**)



- Variables X are separated from Y given Z in graph H, sep_H(X;Y|Z), if there is no active path between any X∈X and any Y∈Y given Z
- The **global Markov assumption** for a Markov network *H* is

10-708 - ©Carlos Guestrin 2006

31

The BN Representation Theorem

If conditional independencies in BN are subset of conditional independencies in P

Obtain

Joint probability distribution:

$$P(X_1,\ldots,X_n) = \prod_{i=1}^n P(X_i \mid \mathbf{Pa}_{X_i})$$

Important because:

Independencies are sufficient to obtain BN structure G

If joint probability distribution:

 $P(X_1,\ldots,X_n) = \prod_{i=1}^n P(X_i \mid \mathbf{Pa}_{X_i})$



Then conditional independencies in BN are subset of conditional independencies in P

Important because:

Read independencies of P from BN structure G

Markov networks representation Theorem 1

If joint probability distribution *P*:

Then

H is an I-map for P

$$P(X_1,\ldots,X_n) = \frac{1}{Z} \prod_{i=1}^m \pi_i(\mathbf{D}_i)$$

■ If you can write distribution as a normalized product of factors ⇒ Can read independencies from graph

10-708 - ©Carlos Guestrin 2006

33

What about the other direction for Markov networks?

If H is an I-map for P

Then

joint probability distribution P:

$$P(X_1,\ldots,X_n) = \frac{1}{Z} \prod_{i=1}^m \pi_i(\mathbf{D}_i)$$

- $\begin{array}{lll} \bullet & \text{Counter-example: } X_1, \ldots, X_4 \text{ are binary, and only eight assignments} \\ & \text{have positive probability:} & \tiny{ \begin{pmatrix} 0,0,0,0 \\ (0,0,0,1) \end{pmatrix} & \tiny{ \begin{pmatrix} 1,1,0,0 \\ (0,0,1,1) \end{pmatrix} & \tiny{ \begin{pmatrix} 1,1,0,0 \\ (0,1,1,1) \end{pmatrix} & \tiny{ \begin{pmatrix} 1,1,1,0 \\ (1,1,1,1) \end{pmatrix} } \\ \end{array}$
- For example, $X_1 \perp X_3 | X_2, X_4$:
- But distribution doesn't factorize!!!

0-708 – ©Carlos Guestrin 2006

Markov networks representation Theorem 2 (Hammersley-Clifford Theorem)

If *H* is an I-map for *P* and *P* is a positive distribution

Then

joint probability distribution *P*:

$$P(X_1,\ldots,X_n) = \frac{1}{Z} \prod_{i=1}^m \pi_i(\mathbf{D}_i)$$

■ Positive distribution and independencies ⇒ P factorizes over graph

10-708 - @Carlos Guestrin 2006

25

Representation Theorem for Markov Networks

If joint probability distribution P:

 $P(X_1,\ldots,X_n) = \frac{1}{Z} \prod_{i=1}^m \pi_i(\mathbf{D}_i)$

Then

H is an I-map for P

If H is an I-map for P and P is a positive distribution

Then

joint probability distribution *P*:

 $P(X_1,\ldots,X_n) = \frac{1}{Z} \prod_{i=1}^m \pi_i(\mathbf{D}_i)$

10-708 - ©Carlos Guestrin 2006

Completeness of separation in Markov networks

■ Theorem: Completeness of separation

- \square For "almost all" distributions that P factorize over Markov network H, we have that I(H) = I(P)
- "almost all" distributions: except for a set of measure zero of parameterizations of the Potentials (assuming no finite set of parameterizations has positive measure)
- Analogous to BNs

10-708 - ©Carlos Guestrin 2006

37

What are the "local" independence assumptions for a Markov network?

■ In a BN *G*:

- local Markov assumption: variable independent of non-descendants given parents
- □ d-separation defines global independence
- □ Soundness: For all distributions:

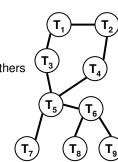
In a Markov net H:

- □ **Separation** defines global independencies
- □ What are the notions of local independencies?

10-708 - ©Carlos Guestrin 2006

Local independence assumptions for a Markov network

- Separation defines global independencies
- Pairwise Markov Independence:
 - □ Pairs of non-adjacent variables are independent given all others



- Markov Blanket:
 - □ Variable independent of rest given its neighbors

10-708 = @Carlos Guestrin 2006

39

Equivalence of independencies in Markov networks

- - **Soundness Theorem**: For all positive distributions *P*, the following three statements are equivalent:
 - ☐ P entails the global Markov assumptions
 - $\ \square$ P entails the pairwise Markov assumptions
 - ☐ P entails the local Markov assumptions (Markov blanket)

0-708 – ©Carlos Guestrin 2006

Minimal I-maps and Markov Networks



- A fully connected graph is an I-map
- Remember minimal I-maps?
 - \square A "simplest" I-map \rightarrow Deleting an edge makes it no longer an I-map
- In a BN, there is no unique minimal I-map
- Theorem: In a Markov network, minimal I-map is unique!!
- Many ways to find minimal I-map, e.g.,
 - □ Take pairwise Markov assumption:
 - ☐ If P doesn't entail it, add edge:

10-708 - ©Carlos Guestrin 2006

41

How about a perfect map?



- Remember perfect maps?
 - \Box independencies in the graph are exactly the same as those in P
- For BNs, doesn't always exist
 - □ counter example: Swinging Couples
- How about for Markov networks?

I0-708 – ©Carlos Guestrin 2006

Unifying properties of BNs and MNs



BNs:

- ☐ give you: V-structures, CPTs are conditional probabilities, can directly compute probability of full instantiation
- but: require acyclicity, and thus no perfect map for swinging couples

MNs:

- □ give you: cycles, and perfect maps for swinging couples
- but: don't have V-structures, cannot interpret potentials as probabilities, requires partition function

Remember PDAGS???

- □ skeleton + immoralities
- □ provides a (somewhat) unified representation
- see book for details

0-708 - @Carlos Guestrin 2006

43

What you need to know so far about Markov networks



- Markov network representation:
 - □ undirected graph
 - □ potentials over cliques (or sub-cliques)
 - □ normalize to obtain probabilities
 - need partition function

Representation Theorem for Markov networks

- □ if P factorizes, then it's an I-map
- □ if P is an I-map, only factorizes for positive distributions
- Independence in Markov nets:
 - □ active paths and separation
 - □ pairwise Markov and Markov blanket assumptions
 - □ equivalence for positive distributions
- Minimal I-maps in MNs are unique
- Perfect maps don't always exist

10-708 - ©Carlos Guestrin 2006