

Inference in graphical models

16-732
Oct 13, 99

Graphical model is joint distribution over domain V .

$$V = U \cup H \cup E$$

$\begin{cases} \rightarrow \text{observed variables (evidence)} \\ \rightarrow \text{hidden variables (unobserved, don't care)} \\ \rightarrow \text{query variables (visible variables) (unobserved, interesting)} \end{cases}$

Inference = answering queries about U , knowing that $E = e$.

Therefore, if P_V is the graphical model, the distribution

$$P_{U|E} = \frac{\int_{\Omega(H)} P_{U,H,E}}{\int_{\Omega(H \cup U)} P_{U,H,E}}$$

is central to the inference methods that we study.

Types of queries: $P_{U|E=e} = ?$

$$P_{x|E=e} = ? \quad x \in U$$

$$\operatorname{argmax}_{\Omega(U)} P_{U|E=e}$$

$$\operatorname{argmax}_{\Omega(x)} P_{x|E=e} \quad x \in U$$

" $E=e$ " is called evidence. This type of evidence is called

categorical. Sometimes evidence is given in the form

of a likelihood, i.e. \tilde{P}_E distribution over $\Omega(E)$.

$$\text{categorical evidence} \Leftrightarrow \tilde{P}_E = \begin{cases} 1 & \text{for } E=e \\ 0 & \text{otherwise} \end{cases}$$

Brute force computation of $P_{U|E}$ (binary variables) 16-732
Oct 13, 99

$$P_{U|E}(u) = \frac{\sum_{h \in \Omega(H)} P_{U,H,E}(u, h, e)}$$

$$\sum_{h \in \Omega(H)} \sum_{u \in \Omega(U)} P_{U,H,E}(u, h, e)$$

$$|\Omega(H)| = 2^{|H|}$$

$$|\Omega(U)| = 2^{|U|}$$

$\Rightarrow O(2^{|H|+|U|})$ operations

exponential in the number of variables that are marginalized over

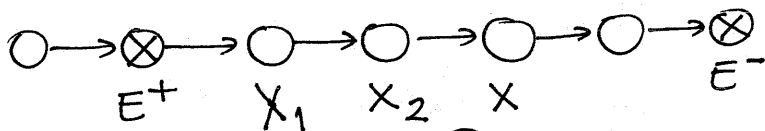
Can we do better by exploiting the structure of graphical models? Yes, if structure is sparse (small no. of parents per variable, small clique size, etc)

Exploiting structure = exploiting independence

Pearl's algorithm

- Bayes nets, singly connected
- computes $P_{X|E}$ for each $X \in V \setminus E$
- linear time (in $|V|$)

Pearl's algorithm in a chain

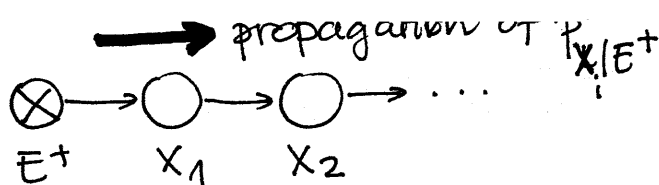


First idea:
$$P_{X|E^+E^-} = \frac{P_{X|E^+E^-}}{P_{E^+E^-}} = \frac{P_{E^+} \cdot P_{X|E^+} \cdot P_{E^-|X,E^+}}{P_{E^+E^-}}$$

$$\propto P_{X|E^+} \cdot P_{E^-|X}$$

second idea: recursion

Computing $P_{X|E^+}$



16-732
Oct 13, 94

$P_{X_1|E^+}$ known

$$P_{X_2|E^+}(x_2|e^+) = \sum_{x_1} P_{X_1|E^+}(x_1|e^+) \cdot P_{X_2|X_1}(x_2|x_1)$$

... and so on, recursively until X

Note that: - computation is recursive, propagating $P_{X_i|E^+}$ down the chain

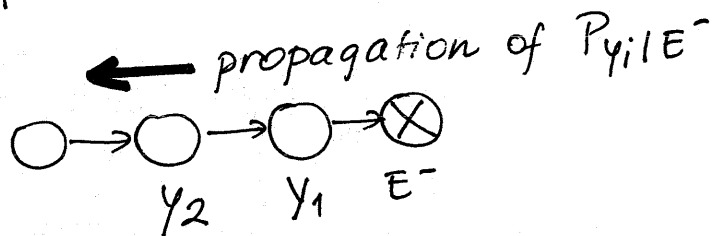
- every stage involves only 2 variables: X_i and its parent, $P_{X_i|X_{i-1}}$ and $P_{X_{i-1}|E^+}$ (the previous "message") \Rightarrow **LOCALITY**

- total nr. of steps $\leq |V| = n$

- computations / step $|\Omega(X_i)| \cdot |\Omega(X_{i-1})|$

if $|\Omega(x)| \leq k \Rightarrow$ total computation $\sim nk^2$

Computing $P_{E^-|X}$



$P_{E^-|Y_1}$ known

$$P_{E^-|Y_2} = \sum_{y_1} P_{E^-|Y_1} \cdot P_{Y_1|Y_2}$$

... and so on

Algorithm

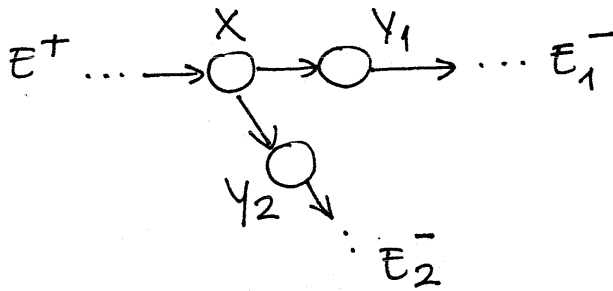
1. Propagate $P_{X|E^+}$ down the chain

2. Propagate $P_{E^-|X}$ up the chain

3. Compute $P_{X|E^+E^-}$ by normalizing $P_{X|E^+} \cdot P_{E^-|X}$

Pearl's algorithm in a tree

16-732
Oct 13, 99



up-propagation E^- is now $\{E_1^-, E_2^-\}$
but $E_1^- \perp E_2^- \mid X$

$$P_{E_1^- E_2^- \mid X} = P_{E_1^- \mid X} \cdot P_{E_2^- \mid X}$$

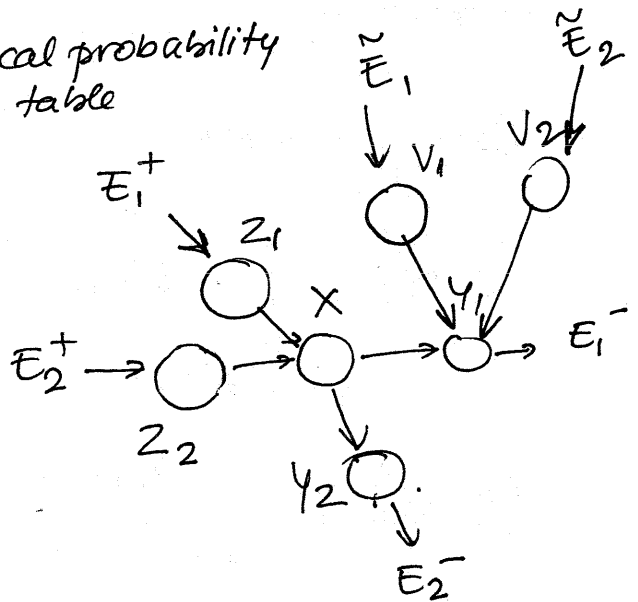
down propagation assume we have $P_{X \mid E^+}$
what is P_{Y_1} "evidence above y_1 " ?

$$P_{Y_1 \mid E_1^+, E_2^-} = \sum_x P_{X \mid E_1^+, E_2^-} \cdot P_{Y_1 \mid X}$$

$$\propto \sum_x P_{X \mid E^+} \cdot P_{E_2^- \mid X} \cdot P_{Y_1 \mid X}$$

involves: downward msg. of parent, upward msgs of siblings, local probability table

Pearl's algorithm in polytrees



polytree = singly connected network

(see [Pearl98] chapter 4, section 5)

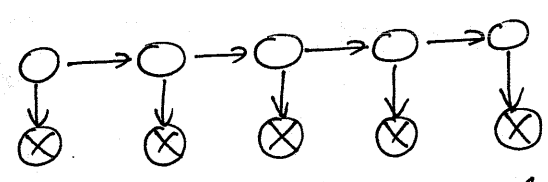
Running time / step $\propto n^{pa(x)}$

Conclusion about Pearl's algorithm

16-732
Oct 13, 96

- computes $P_{X|E}$ for all $X \in V \setminus E$
- running time :- linear in $n = |V|$
 - polynomial in $k = \max |\Omega(x)|$ if $\max \# \text{ parents bounded}$
 - exponential in $\max \# \text{ parents of } X \in V$
- involves only local computations
 - local variables
 - local prob. tables (parameters)
 - messages from parents, children
- makes use of independence relationships
- works exactly on singly connected Bayes nets (polytrees)

Special cases



Hidden Markov Model

Pearl's alg \Leftrightarrow Forward-Backward Alg.
(can do max propagation as well \Leftrightarrow Viterbi)

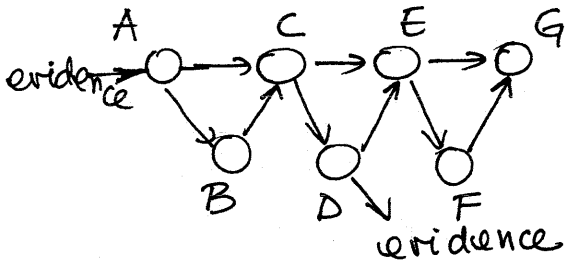
HMM + continuous variables + gaussian distribution + linear dependencies
 \Rightarrow Linear system (discrete time)

\Downarrow
Pearl's algorithm \Leftrightarrow [batch version of] Kalman filter

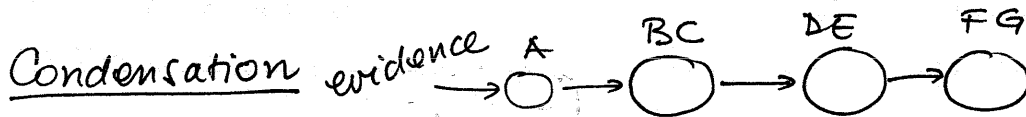
Pearl's alg. itself is a special case of the junction tree alg that will be discussed later on.

What to do on Bayes nets with loops?

16-732
Oct 13, 99



- Pearl's algorithm — approximation!
- conditioning = breaking the loops
- condense variables
- convert to junction tree —
- general method

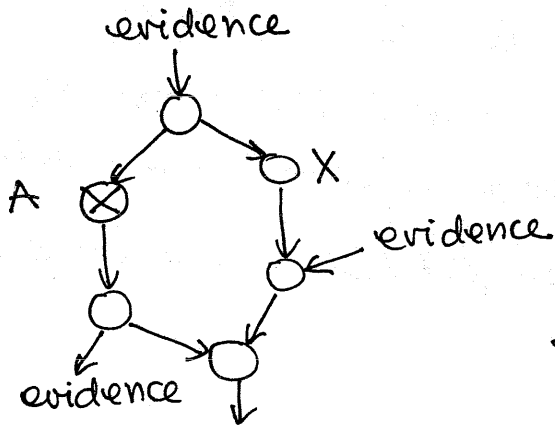


computation cost \sim exponential in the size of largest ^{set of} condensed variables

Junction tree alg is a special case of condensation + propagation

Conditioning

assume A known \Rightarrow blocks a path and breaks cycle



$$P_{X|evidence} = P_{X|A=0, evid} \cdot P(A=0|evidence) + P_{X|A=1, evid} \cdot P(A=1|evidence)$$

$$P(A=0|evidence) \propto P(evidence|A=0) \cdot P_A$$

computation cost \sim exponential in the number of instantiated variables

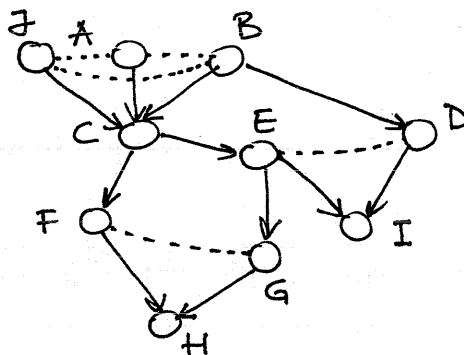
Converting a Bayes net into a junction tree

16-732
Oct 13, 99

- by adding edges \Rightarrow hides some independence relationships
(\Rightarrow they won't be exploited in the inference procedure)

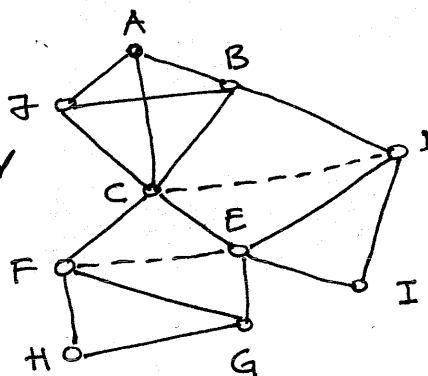
1. Moralization

- "marry" parents of common child
- then drop direction of edges



2. Triangulation

- add edges until graph triangulated
- triangulation is not unique!



3. Junction tree construction

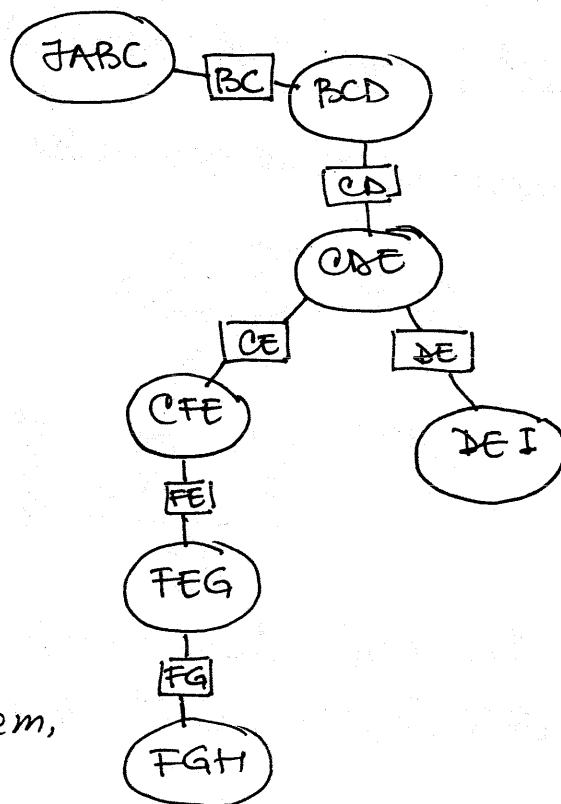
- junction tree may not be unique (same cliques, same separators but different topologies)

Total "state space" size =

$$= \sum_{C \in \text{cliques}} |\Omega(C)|$$

indicator of total memory, propagation time

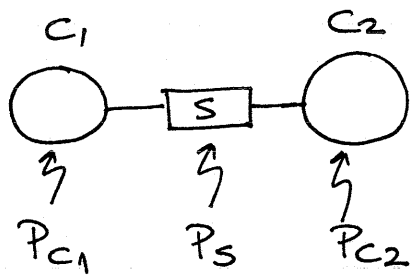
Wanted triangulations that give small total state space. Hard problem, only heuristics available.



The Junction tree algorithm

16-732
Oct 13, 1996

Passing messages between two cliques



before propagation:

$$P_S = \sum_{\Omega(C_1 \setminus S)} P_{C_1} = \sum_{\Omega(C_2 \setminus S)} P_{C_2}$$

we say that junction tree is "balanced"

new info: $P_{C_1} \leftarrow P_{C_1}^*$

Propagating this new info to C_2 is called absorption.

$$1. P_S^* = \sum_{\Omega(C_1 \setminus S)} P_{C_1}^*$$

new separator potential

$$2. P_{C_2}^* = P_{C_2} \cdot \frac{P_S^*}{P_S}$$

replaces P_{C_2} ; then P_S^* replaces P_S

Why does this operation make sense?

First, note that absorption does not change the joint distribution.

Supposing there are only two cliques, i.e. $V = C_1 \cup C_2$:

$$P_V = \frac{P_{C_1}^* \cdot P_{C_2}}{P_S} \quad \text{before absorption}$$

$$P_V^{\text{new}} = \frac{P_{C_1}^* \cdot \left(P_{C_2} \cdot \frac{P_S^*}{P_S} \right)}{P_S^*} = P_V \quad \text{after absorption}$$

Thus, absorption distributes info that is already in the distribution. But how does this info get incorporated?

Second Assume that $P_V = \frac{P_{C_1} \cdot P_{C_2}}{P_S}$ and

16-732
Oct 13, 199

you observe some [subset of] variable[s] $A \in C_1 \setminus C_2$.

Hence $A=a$. Before the observation we had

$$P_A = \sum_{\Omega(C_1 \setminus A)} P_{C_1} = \text{the marginal of } A \text{ computed in } P_{C_1}$$

After the observation, we have a new " P_A ":

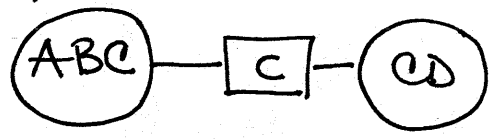
$$P_A^* = \delta_a = \begin{cases} 1 & A=a \\ 0 & \text{otherwise} \end{cases}$$

We incorporate this into C_1 by the law of conditional probability

(let's denote $C_1 \setminus A = B$)

$$P_{C_1}^* = P_{C_1 | \text{observation}}(b, a') = \begin{cases} \frac{P_{C_1}(b, a)}{P_A} & \text{if } a' = a \\ 0 & \text{if } a' \neq a \end{cases}$$

Example:



A, B, C, D binary variables

		A	
		0	1
P _{ABC} →	00	0.2	0.1
	01	0.15	0.05
	11	0.2	0.2
P _C →	10	0.05	0.05

		A	
		0	1
P _C =		0.4	0.6

		C	
		0	1
P _{CD} =	D = 0	0.1	0.5
	D = 1	0.3	0.1

Observed $A=0$

16-732
Oct 13, 19

$$P_A(a=0) = 0.2 + 0.15 + 0.2 + 0.05 = 0.6$$

$$P_{ABC}^* = \frac{P_{ABC} \cdot \delta_{a,0}}{P_A} = \frac{1}{0.6}$$

0.2	0
0.15	0
0.2	0
0.05	0

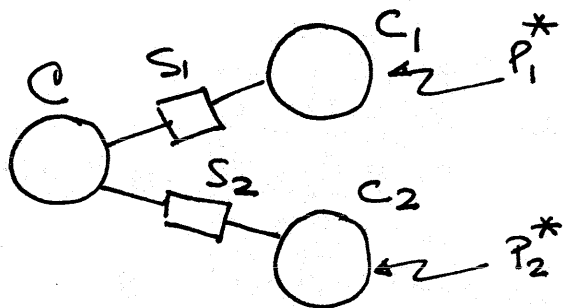
$$P_C^* = \sum_{a,b} P_{ABC}^*(a,b,c) = \begin{array}{|c|c|} \hline \frac{0.25}{0.6} & \frac{0.35}{0.6} \\ \hline \end{array}$$

$$P_{CD}^* = P_{CD} \cdot \frac{P_C^*}{P_C} = \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 0 & 0.104 & 0.486 \\ \hline 1 & 0.312 & 0.097 \\ \hline \end{array}$$

for ex: $P_{CD}(0,1) = 0.5$

$$P_{CD}^*(0,1) = P_{CD}(0,1) \cdot \frac{P_C^*(0)}{P_C(0)} = 0.5 \cdot \frac{\frac{0.25}{0.6}}{0.4} = 0.52$$

How to absorb from several cliques?



$$1. P_{S_i}^* = \sum_{R(C_i | S_i)} P_{C_i}^* \quad i=1,2,\dots$$

$$2. P_C^* = P_C \cdot \frac{P_{S_1}^*}{P_{S_1}} \cdot \frac{P_{S_2}^*}{P_{S_2}} \dots$$

This is justified by the requirement that P_V remain constant after absorption.

The general case

Arbitrary junction tree

$$P_V = \frac{\prod P_{C_i}}{\prod P_{S_j}}$$

16-732
Oct 13, 1996

Arbitrary categorical evidence

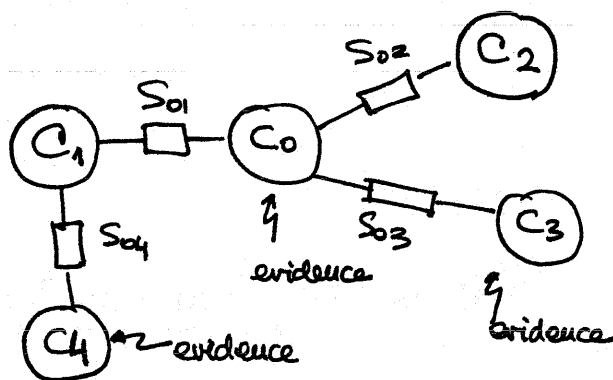
$$X_k = x_k^*$$

$$k = 1, \dots, n_e$$

1. Incorporate evidence:

for each k find C_i clique such that $X_k \in C_i$

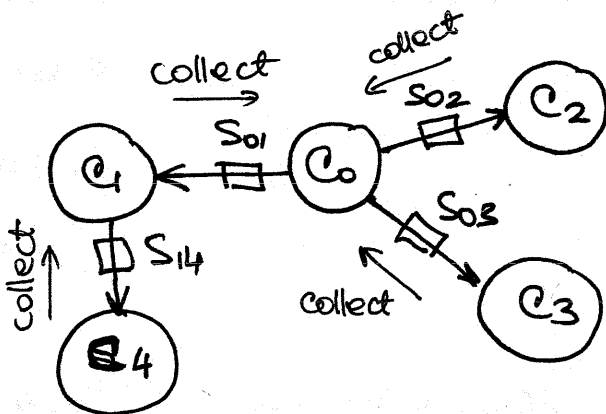
$$P_{C_i}^* = \frac{P_{C_i} \cdot \delta_{X_k, x_k^*}}{P_{X_k}(x_k^*)}$$



2. "Direct" the junction tree.

- choose [any clique] C_0 as root

- direct all edges away from root



3. Collect evidence

starting from root do recursively:

- call "collect evidence" in all children
- then absorb from all children

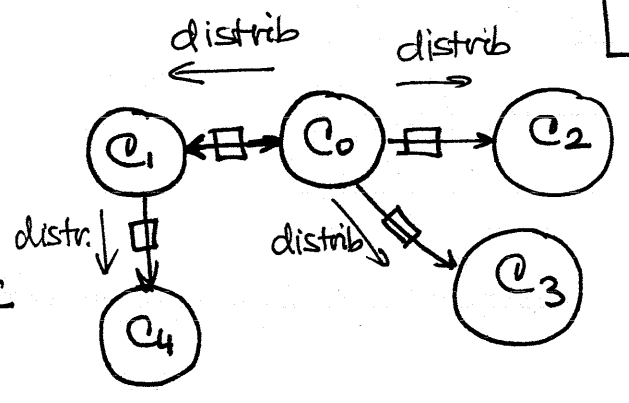
(this results in a series of absorptions from leaves towards the root. For ex: C_1 absorbs from C_4

$$C_0 \text{ --- } C_1, C_2, C_3$$

4. "Distribute evidence"

starting from root do recursively
for each clique C :

- each child of C absorbs from C
- each child calls "distribute evidence"



For ex, here : C_1, C_2, C_3 absorb from C_0
 C_4 absorbs from C_1

END

This is the Junction Tree Algorithm.

- It ensures that :
- the tree is balanced
 - every P_C^*, P_S^* is normalized
 - every P_C^* is the marginal of C in P_V^*
 - P_V^* represents $\text{Prob}[V | \text{evidence}]$

Non-categorical evidence :

• info from external source tells us

$$P_{X_k}^* \quad (\neq P_{X_k})$$

AND

$$(\neq \delta_{X_k, x_k^*})$$

• then $P_{C_i}^* = \frac{P_{C_i} \cdot P_{X_k}^*}{P_{X_k}}$ in (1.) and the other
steps remain unchanged

Computational complexity of the junction tree algorithm

16-732
Oct 13, 19

Absorption $C_1 \rightarrow C_2$ takes $\sim |\Omega(C_1)|$ additions
 $|\Omega(C_2)|$ multiplications

Collect evidence: absorptions children \rightarrow parent $\parallel \Rightarrow$
Distribute evidence: \longleftarrow parent \rightarrow child

whole alg $\approx 2 \cdot \underbrace{\sum_{C \in \text{cliques}} |\Omega(C)|}_{2 \cdot T}$ additions and multiplications

T is called "total state space"; it's an approximate measure of the cost of inference in the junction tree.

- Note that:
- $|\Omega(C)|$ is exponential in $|C|$
 - therefore $T \ll |\Omega(V)|$ if tree is sparse
 \uparrow cost of "brute force" inference
 - it is important to minimize T during triangulation. (but this is hard!)
 - a good approximation to T is $\max_C |\Omega(C)|$
(minimizing this is also NP-hard)

The JT algorithm • is exact and general

- any exact, general algorithm for inference has to perform a [possibly hidden] triangulation; hence, it is essentially no better than the JT algorithm
- unfortunately, for many practical problems the JT is intractable: image segmentation, max likelihood decoding, medical diagnosis in QMR-DT

Approximate inference methods

16-732
Oct 13, 1996

- Pearl's (polytree) algorithm ← for Bayes nets
 - comes with no guarantees
 - seems to work in practice
- Markov Chain Monte Carlo methods (MCMC)
i.e. sampling ← for Markov fields mainly
- variational methods ← for Bayes nets, Markov fields
approximating $P_V / \text{evidence}$

Graphical Probability Models (aka Belief Nets)

16-732
Oct 6, 94

Pearl (1988):

"probability theory is unique in its ability to process context-sensitive beliefs, and what makes the process computationally feasible is that the information needed for specifying context dependencies is represented by graphs and manipulated by local propagation".

Graphical model = graphical representation of (conditional) independence relationships in a joint distribution
= the distribution itself

graphical model / structure (a graph)
graphical model \ parametrization (depends on graph, parameters are "local")

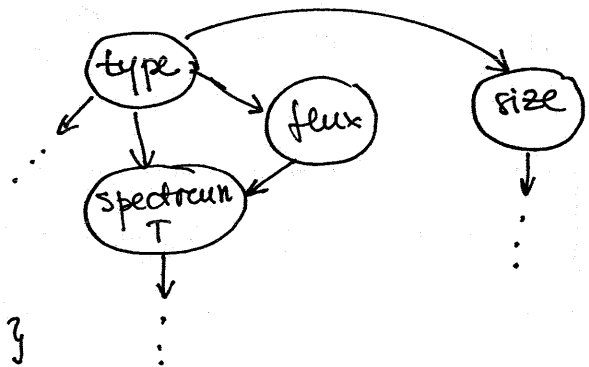
graph = (V , E)
 ↓ ↓
 nodes edges
 variables "dependencies"

More precisely: a missing edge represents an independence relationship.

But what is "independence"?

Example: "Redshift" model

V = { type, flux, size, ... }
E = { type → flux, type → size, ... }



Probabilistic independence

16-732
Oct 6, 99

$$A \perp B \iff P_{AB} = P_A \cdot P_B$$

"A independent B"

(marginal formulation)

Example 1. A, B binary variables, $A \perp B$

		A		P_B ↓
		0	1	
P_{AB}	B	0	1	
	0	0.12	0.28	0.4
	1	0.18	0.42	0.6
P_A		0.3	0.7	

Example 2. A, B real, jointly gaussian

$$[A \ B] \sim N([\mu_A \ \mu_B], \Sigma_{AB})$$

$$P_{AB}(a, b) = \frac{1}{|\Sigma|^{1/2} \cdot 2\pi} \exp \left\{ -\frac{1}{2} (a - \mu_A \ b - \mu_B) \Sigma_{AB}^{-1} \begin{pmatrix} a - \mu_A \\ b - \mu_B \end{pmatrix} \right\}$$

$$A \perp B \iff \Sigma_{AB} = \begin{bmatrix} \sigma_A^2 & 0 \\ 0 & \sigma_B^2 \end{bmatrix} \text{ (diagonal)}$$

Ex 1. Prove this!

Ex 2. $A \perp B \stackrel{?}{\implies} A \perp B, C$
 $A \perp C \stackrel{?}{\impliedby}$

Is any of the implications \implies true?

Prove or give counterexample!

$$A \perp B \iff P_{A|B} = P_A \iff$$

↑
(conditional formulation)

knowing B brings no additional information about A

Conditional independence

16-732
Oct 6, 9

$$A \perp B \mid C \iff P_{A|B|C} = P_{A|C} \cdot P_{B|C} \iff P_{A|BC} = P_{A|C}$$

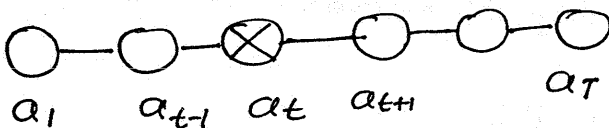
"once C is known, knowing B brings no additional information about A"

→ has to hold for all values of A, B, C

Representing independences in graphs

Independence (in joint distribution) \longleftrightarrow Separation (in graph)

Markov chain

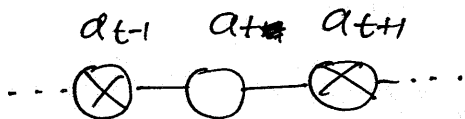


$$i < t < j$$

$$a_i \perp a_j \mid a_t$$

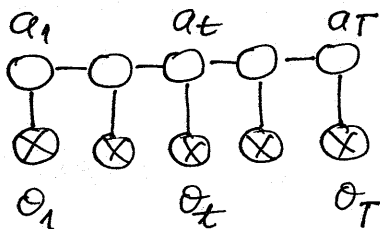
Knowing a state makes future independent from past.

Knowing $a_{t+1}, a_{t-1} \Rightarrow$ no other state can give additional info about a_t



$\Rightarrow a_{t-1}, a_{t+1}$ separate a_t from the other variables

Variant Hidden Markov Model



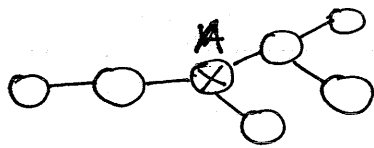
Separation properties allow for efficient algorithms for HMMs (forward-backward, Viterbi, Kalman filter, Baum-Welch)

There are T "hidden" variables. $\Rightarrow 2^T$ hidden states

16-732
Oct 6, 1996

But, forward-backw, etc, ... are $O(T)$ (i.e. linear!)

Tree



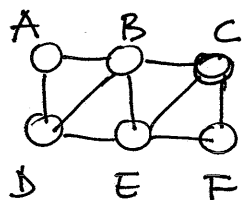
graph with no cycles
(or singly connected)

Property - every variable A breaks the tree into 2 or more subtrees

- between two variables there is exactly 1 path

Trees are the only structure for which learning can be done efficiently

Markov field



Separation

• $V_1 \perp V_2 \mid V_3 \iff V_3$ "blocks" all the paths from V_1 to V_2

• variable \perp everything | neighbors
else

(i.e. all info about a variable X that can be obtained from $V \setminus \{X\}$ is contained in the neighbors of X)

Examples:

$$A, D \perp C, F \mid B, E$$

$$B \perp F \mid C, E$$

$$B \not\perp F \mid D, E$$

$$A \perp C, E, F \mid \underbrace{B, D}_{\text{neighbors of } A}$$

Parametrization

Clique = fully connected subgraph of a graph

Ex: • all edges are size 2 cliques

• ABD, BDE, BCE, CEF are size 3 cliques. They are maximal.

• $BCDE$ is not a clique

Denote by \mathcal{C}_G the set of cliques of a graph G .
Maximal

\mathcal{P}_V Markov Field
over graph G

\Leftrightarrow

$$\mathcal{P}_V(x) = \prod_{c \in \mathcal{C}_g} \psi_c(x_c)$$

16-732
Oct 6, 1999

Ex: $\mathcal{P}_{ABCDEF} = \psi_{ABD}(abd) \psi_{BDE}(bde) \psi_{BCE}(bce) \psi_{CEF}(cef)$

↓
"local" function/parametrization - depends only on the variables ABD

How much memory is saved?

- assume A, B, \dots, F binary variables
- assume ψ are [probability] tables

\mathcal{P}_{ABCDEF} = prob. table with 2^6 entries, $2^6 - 1$ free parameters

ψ_{ABD} = prob table with 2^3 entries, $2^3 - 1$

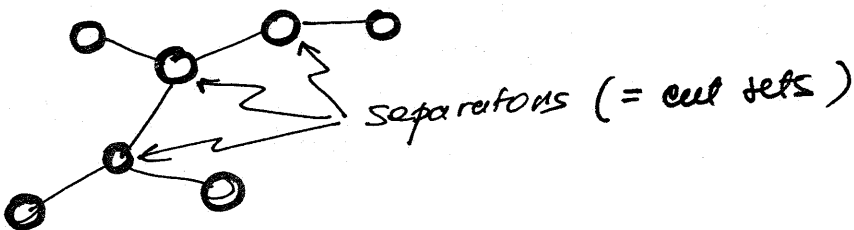
saving: $2^6 - 1 - 4(2^3 - 1) = 35$ out of 63 (more than 50%)

Ex: Tree over n variables. (binary)

$T \equiv \mathcal{P}_V = \prod_{AB \in E} \psi_{AB}$ Memory: $3(n-1)$

For a tree, ψ_{AB} can be related to the corresponding marginals \mathcal{P}_{AB}

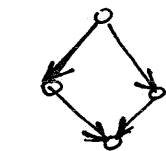
$$T = \frac{\prod_{AB \in E} \mathcal{P}_{AB}(a, b)}{\prod_{A \in V} \mathcal{P}_A^{\text{deg } A - 1}(a)} = \frac{\prod \mathcal{P}_{\text{clique}}}{\prod \mathcal{P}_{\text{separator}}}$$



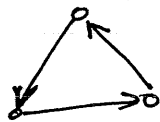
Bayes Nets

Oct 6, 1996
16-732

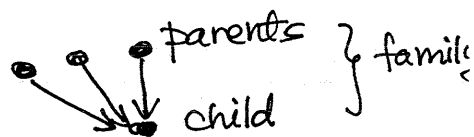
DAG (Directed Acyclic Graph) = directed graph with no directed cycles



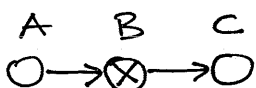
a DAG



not a DAG

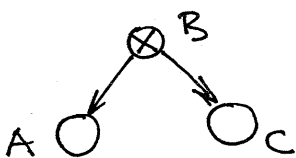


Separation in a DAG (d-separation)

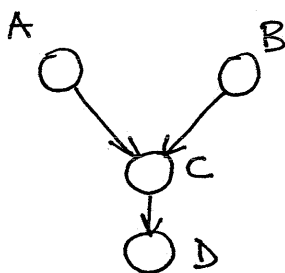


$$A \perp C \mid B$$

(path $A \rightarrow B \rightarrow C$ blocked)



$$A \perp C \mid B, \text{ but } A \not\perp C$$



$$A \perp B \text{ but } A \not\perp B \mid C$$

$$A \not\perp B \mid D$$

(path $A \rightarrow B$ blocked by C or D)

In general $A \perp B \mid C$ in a DAG iff all paths btw. A and B are blocked when C is instantiated.

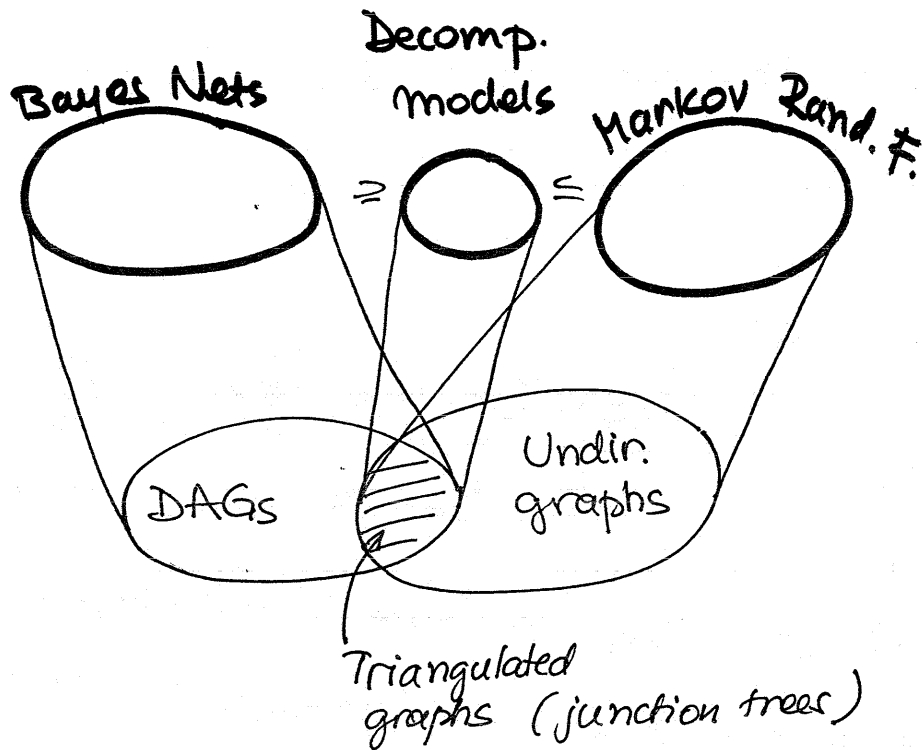
Understanding separation is important for inference: if two variables are separated we can ignore one when reasoning about the other.

Parametrization

$$P_V = \prod_{A \in V} P_{A \mid \text{pa}(A)}$$

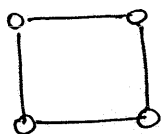
local fctns (depend only on A, pa(A))

$\text{pa}(A)$ = parents of variable A

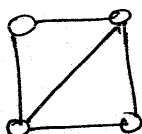


Triangulated graphs / decomposable models

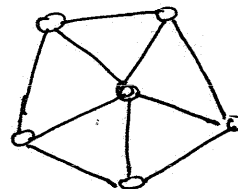
Graph G triangulated (or chordal) \Leftrightarrow every cycle ≥ 4 in G has a chord



not triangulated

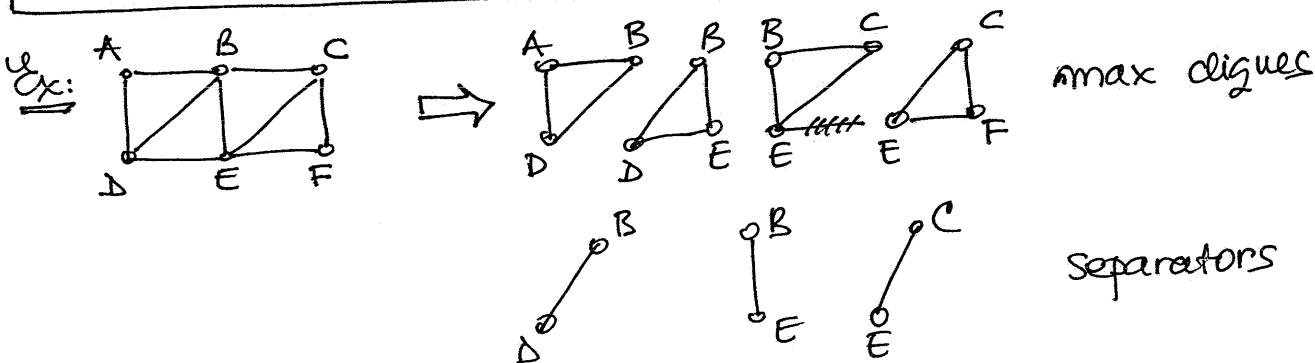


triangulated



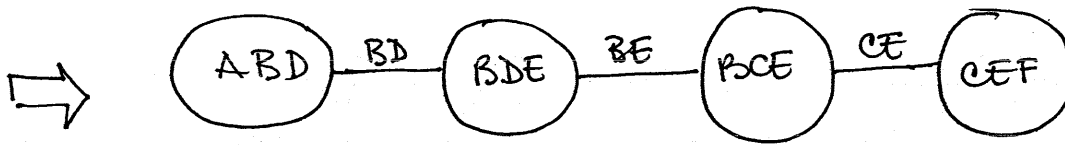
NOT triangulated

Triangulated graphs can be represented as a tree of max cliques.



Junction tree

16-732
Oct 6, 199



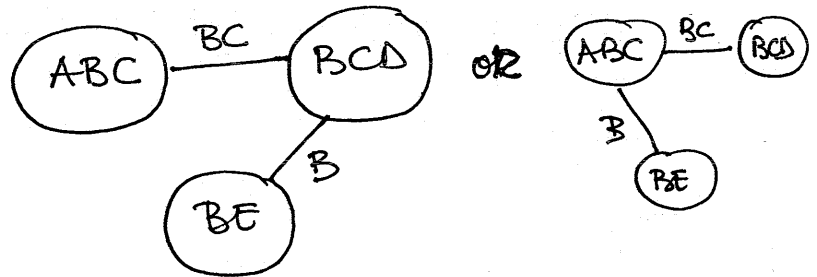
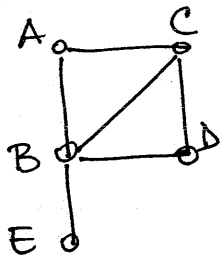
⇒ parametrization:

$$P_{ABCDEF} = \frac{P_{ABD} \cdot P_{BDE} \cdot P_{BCE} \cdot P_{CEF}}{P_{BD} \cdot P_{BE} \cdot P_{CE}} = \frac{\prod_C P_{cliques}}{\prod_S P_{separators}}$$

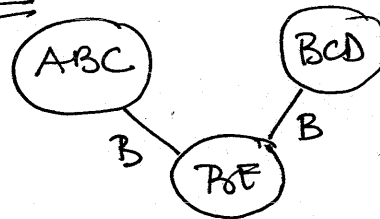
Remarks: • Junction tree structure is not unique

• Junction tree is a max spanning tree w.r.t. separator size

Ex:



but NOT:



The Modified Pearl Algorithm

(After Peot & Schachter)

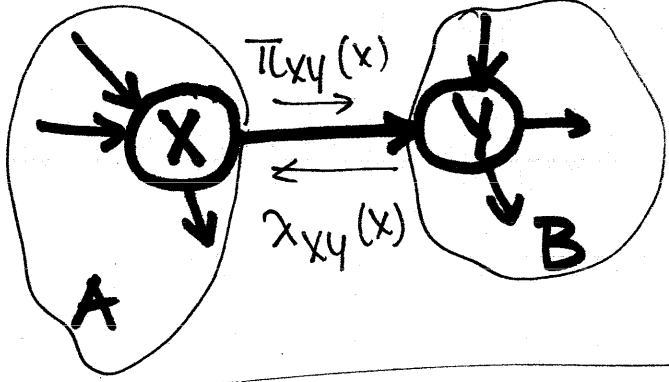
At each node X:

$U = \{U_1 \dots U_m\}$ parents

$Y = \{Y_1 \dots Y_n\}$ children

$P_{X|U}$ conditional prob. table

$\lambda_X(x) =$ evidence for X
 $= \begin{cases} \delta_x & \text{if } X=x \text{ observed} \\ 1 & \text{if } X \text{ unobserved} \end{cases}$



$$\pi_{XY}(x) = P(X=x, e_A)$$

$$\lambda_{XY}(x) = P(e_B | X=x)$$

"MESSAGES"

local data tables: $\pi(x), \lambda(x)$
 $P(x, e^+)$ $P(e^- | X=x)$

Algorithm (Modified Pearl)

Initialize $\pi(x) = \pi_{XY_j}(x) = P(X=x)$ obtained by a propagation without evidence
 $\lambda = 1$

Update at node X:

$$P(x) \propto \pi(x) \lambda(x)$$

$$\pi(x) = \sum_{u \in \Omega_U} P(X=x | U=u) \prod_{k=1}^m \pi_{U_k X}(u_k)$$

$$\lambda(x) = \lambda_X(x) \prod_{i=1}^n \lambda_{XY_i}(x)$$

[Modified Pearl Algorithm]

STAT 593A
Spring 05

To parent U_i

$$\lambda_{U_i X}(u_i) = \sum_x \lambda(x) \sum_{u' \in \Omega_U \setminus U_i} P(x | u', u_i) \prod_{k \neq i} \pi_{U_k X}(u_k)$$

To child Y_j

$$\pi_{X Y_j}(x) = \pi(x) \lambda_X(x) \prod_{k \neq j} \lambda_{X Y_k}(x)$$

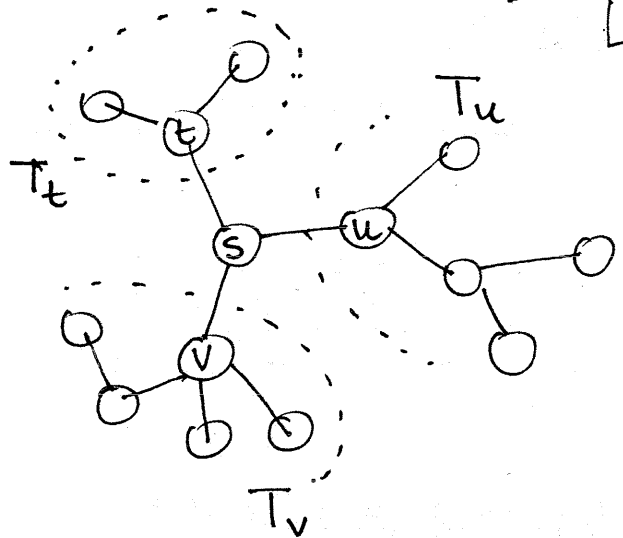
Sum-Product Algorithm

on Trees

(After Wainwright & Jordan)

Tree distribution

$$P_X = \frac{\prod_{st \in E} P_{st}(x_s, x_t)}{\prod_{s \in V} P_s(x_s)^{\text{deg}_s - 1}}$$



"potentials" $\psi_{st} = P_{st}$

$$\psi_s = 1/P_s^{\text{deg}_s - 1}$$

enter evidence by multiplication

$$\Rightarrow P_X(x) = \prod_{s \in V} \psi_s(x_s) \cdot \prod_{st \in E} \psi_{st}(x_s, x_t)$$

Wanted: $\mu_s(x_s) = \sum_{x': x'_s = x_s} P_X(x')$ Marginal of node s

Remarks:

Node s separates tree into s, and disjoint trees T_u, T_v, T_t

$n(s)$

Denote $T_u = (V_u, E_u)$

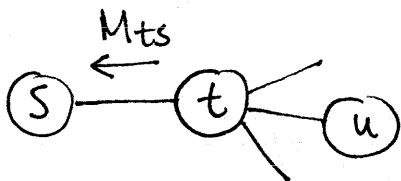
$$\mu_s(x_s) = \sum_{x': x'_s = x_s} \psi_s(x_s) \prod_{u \in n(s)} P_{x_s u}(x'_u)$$

$$= \psi_s(x_s) \cdot \prod_{u \in n(s)} \left[\sum_{x'_u} P_{x_s u}(x'_u) \cdot \psi_{su}(x_s, x'_u) \right]$$

$M_{us}^*(x_s) \leftarrow \dots$ "message" from T_u

Algorithm

Message updating: $M_{ts}(x_s) \leftarrow \sum_{x'_t} \left[\psi_{st}(x_s, x'_t) \psi_t(x'_t) \prod_{u \in \mathcal{N}(t) \setminus s} M_{ut}(x'_t) \right]$



(*)

Message propagation:

- Each node s gets messages from all its neighbors, and sends messages to all neighbors (asynchronous, naive)
- OR: - Choose a root s
 - Collect messages recursively from leaf nodes to s
 - Distribute messages recursively from s toward leaf nodes

Sum-Product: in (*), each message is computed as a product of messages received, followed by a sum.

Max-Product: Replace " \sum " by "max" in (*) \Rightarrow Dynamic programming algorithm to compute the most likely configuration