

Bayesian Networks: Independencies and Inference

*with annotations
from lecture on Thu Nov 3rd*

Scott Davies and Andrew Moore

Note to other teachers and users of these slides. Andrew and Scott would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials>. Comments and corrections gratefully received.

$P(G|B) = P(G \wedge S | B) + P(G \wedge \sim S | B)$
 $P(G|S=i) = \frac{P(B|S=i)P(S=i)}{P(B)}$

$P(X_i=j) \text{ (E)}$

$P(S) = \frac{1}{3}$

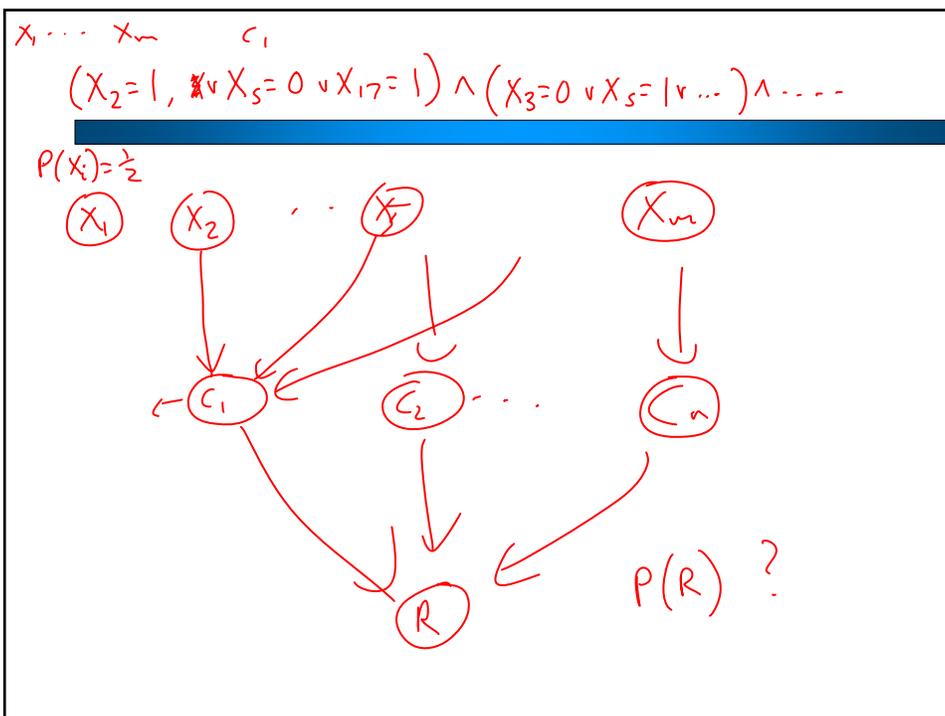
$P(B|S) = \frac{1}{2}$ $P(G|S) = 1$
 $P(B|\sim S) = \frac{1}{4}$ $P(G|\sim S) = \frac{1}{40}$

$P(G=0, B=0, S=0) =$
 $P(G=0|B=0, S=0)P(B=0, S=0) =$
 $\frac{P(G=0|S=0)}{3^2 \cdot 40} \times \frac{3}{4} \times \frac{2}{3}$

Live in Sq Hill?
 $P(G|B) = \frac{P(G \wedge B)}{P(B)} = \frac{P(G \wedge B \wedge S) + P(G \wedge B \wedge \sim S)}{P(B)}$
 $P(\dots) \} \text{35 squid hillers}$

$P(G|B)$ ① Get Joint
 ② Int from Joint

S	B	G	$P(S, B, G)$
0	0	0	...
0	0	1	...
0	1	0	...
0	1	1	...
1	0	0	...
1	0	1	...
1	1	0	...
1	1	1	...



What Independencies does a Bayes Net Model?

- In order for a Bayesian network to model a probability distribution, the following must be true by definition:

Each variable is conditionally independent of all its non-descendants in the graph given the value of all its parents.

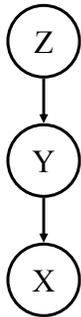
- This implies

$$P(X_1 \dots X_n) = \prod_{i=1}^n P(X_i | \text{parents}(X_i))$$

- But what else does it imply?

What Independencies does a Bayes Net Model?

- Example:



Given Y , does learning the value of Z tell us nothing new about X ?

I.e., is $P(X|Y, Z)$ equal to $P(X | Y)$?

Yes. Since we know the value of all of X 's parents (namely, Y), and Z is not a descendant of X , X is conditionally independent of Z .

Also, since independence is symmetric,
 $P(Z|Y, X) = P(Z|Y)$.

Quick proof that independence is symmetric

- Assume: $P(X|Y, Z) = P(X|Y)$
- Then:

$$P(Z | X, Y) = \frac{P(X, Y | Z)P(Z)}{P(X, Y)} \quad (\text{Bayes's Rule})$$

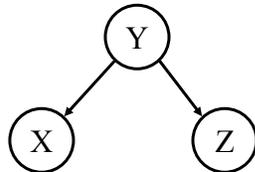
$$= \frac{P(Y | Z)P(X | Y, Z)P(Z)}{P(X | Y)P(Y)} \quad (\text{Chain Rule})$$

$$= \frac{P(Y | Z)P(X | Y)P(Z)}{P(X | Y)P(Y)} \quad (\text{By Assumption})$$

$$= \frac{P(Y | Z)P(Z)}{P(Y)} = P(Z | Y) \quad (\text{Bayes's Rule})$$

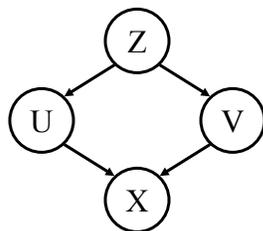
What Independencies does a Bayes Net Model?

- Let $I\langle X, Y, Z \rangle$ represent X and Z being conditionally independent given Y .



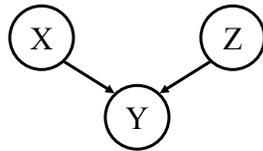
- $I\langle X, Y, Z \rangle$? Yes, just as in previous example: All X 's parents given, and Z is not a descendant.

What Independencies does a Bayes Net Model?



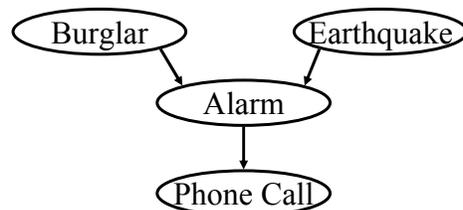
- $I\langle X, \{U\}, Z \rangle$? No.
- $I\langle X, \{U, V\}, Z \rangle$? Yes.
- Maybe $I\langle X, S, Z \rangle$ iff S acts a cutset between X and Z in an undirected version of the graph...?

Things get a little more confusing



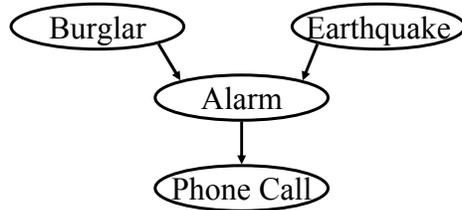
- X has no parents, so we know all its parents' values trivially
- Z is not a descendant of X
- So, $I\langle X, \{\}, Z \rangle$, even though there's a undirected path from X to Z through an unknown variable Y.
- What if we do know the value of Y, though? Or one of its descendants?

The “Burglar Alarm” example



- Your house has a twitchy burglar alarm that is also sometimes triggered by earthquakes.
- Earth arguably doesn't care whether your house is currently being burgled
- While you are on vacation, one of your neighbors calls and tells you your home's burglar alarm is ringing. Uh oh!

Things get a lot more confusing



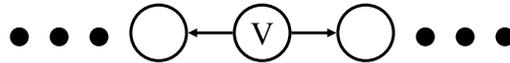
- But now suppose you learn that there was a medium-sized earthquake in your neighborhood. Oh, whew! Probably not a burglar after all.
- Earthquake “explains away” the hypothetical burglar.
- But then it must **not** be the case that $I\langle \text{Burglar}, \{\text{Phone Call}\}, \text{Earthquake} \rangle$, even though $I\langle \text{Burglar}, \{\}, \text{Earthquake} \rangle$!

d-separation to the rescue

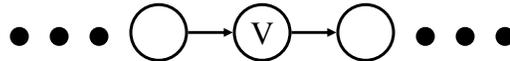
- Fortunately, there is a relatively simple algorithm for determining whether two variables in a Bayesian network are conditionally independent: *d*-separation.
- Definition: X and Z are *d*-separated by a set of evidence variables E iff every undirected path from X to Z is “blocked”, where a path is “blocked” iff one or more of the following conditions is true: ...

A path is “blocked” when...

- There exists a variable V on the path such that
 - it **is** in the evidence set E
 - the arcs putting V in the path are “tail-to-tail”



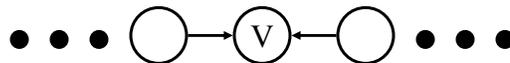
- Or, there exists a variable V on the path such that
 - it **is** in the evidence set E
 - the arcs putting V in the path are “tail-to-head”



- Or, ...

A path is “blocked” when... (the funky case)

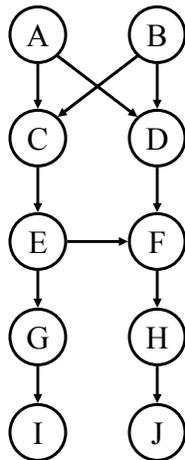
- ... Or, there exists a variable V on the path such that
 - it is **NOT** in the evidence set E
 - **neither are any of its descendants**
 - the arcs putting V on the path are “head-to-head”



d-separation to the rescue, cont'd

- Theorem [Verma & Pearl, 1998]:
 - If a set of evidence variables E d -separates X and Z in a Bayesian network's graph, then $I\langle X, E, Z\rangle$.
- d -separation can be computed in linear time using a depth-first-search-like algorithm.
- Great! We now have a fast algorithm for automatically inferring whether learning the value of one variable might give us any additional hints about some other variable, given what we already know.
 - "Might": Variables may actually be independent when they're not d -separated, depending on the actual probabilities involved

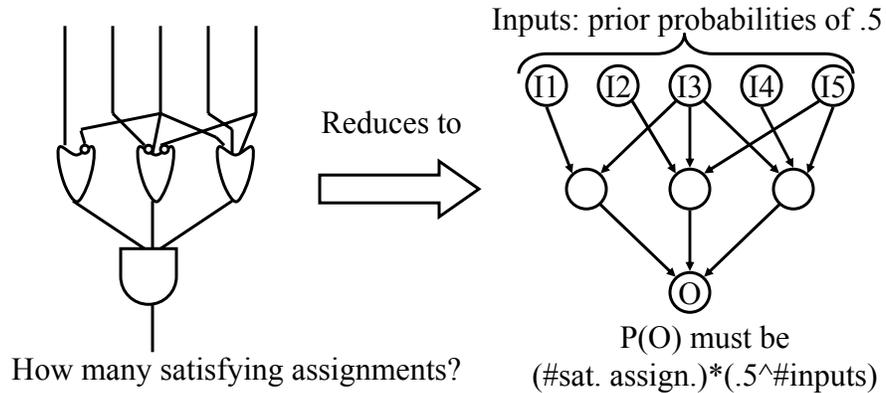
d-separation example



- $I\langle C, \{\}, D\rangle?$
- $I\langle C, \{A\}, D\rangle?$
- $I\langle C, \{A, B\}, D\rangle?$
- $I\langle C, \{A, B, J\}, D\rangle?$
- $I\langle C, \{A, B, E, J\}, D\rangle?$

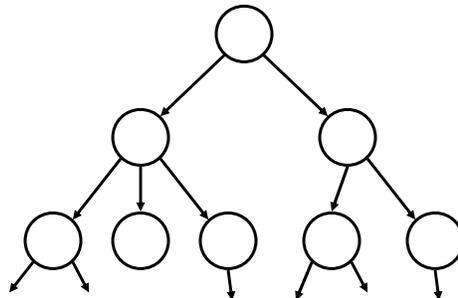
Bayesian Network Inference

- Inference: calculating $P(X|Y)$ for some variables or sets of variables X and Y .
- Inference in Bayesian networks is #P-hard!



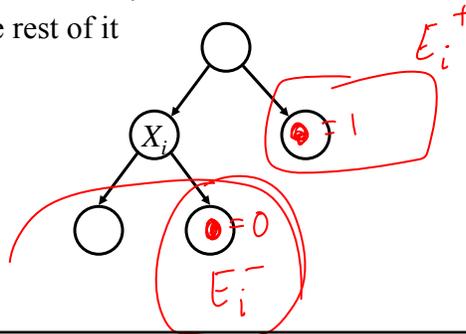
Bayesian Network Inference

- **But...** inference is still tractable in some cases.
- Let's look a special class of networks: *trees* / *forests* in which each node has at most one parent.



Decomposing the probabilities

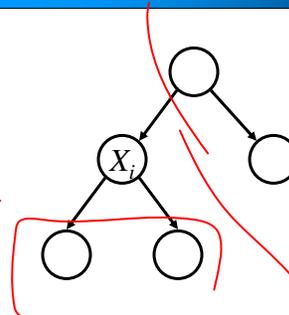
- Suppose we want $P(X_i | E)$ where E is some set of evidence variables.
- Let's split E into two parts:
 - E_i^- is the part consisting of assignments to variables in the subtree rooted at X_i
 - E_i^+ is the rest of it



Decomposing the probabilities, cont'd

$$P(X_i | E) = P(X_i | E_i^-, E_i^+)$$

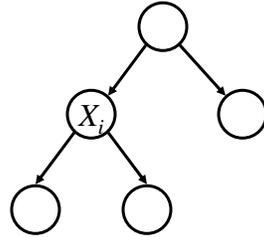
$$= \frac{P(E_i^- | X_i, E_i^+) P(X_i | E_i^+)}{P(E_i^- | E_i^+)}$$



$$= \propto P(E_i^- | X_i) P(X_i | E_i^+)$$

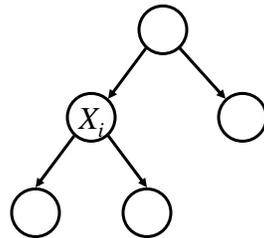
Decomposing the probabilities, cont'd

$$\begin{aligned} P(X_i | E) &= P(X_i | E_i^-, E_i^+) \\ &= \frac{P(E_i^- | X, E_i^+) P(X | E_i^+)}{P(E_i^- | E_i^+)} \end{aligned}$$



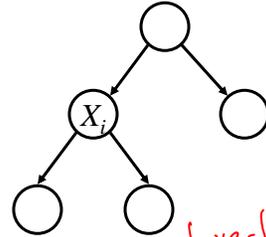
Decomposing the probabilities, cont'd

$$\begin{aligned} P(X_i | E) &= P(X_i | E_i^-, E_i^+) \\ &= \frac{P(E_i^- | X, E_i^+) P(X | E_i^+)}{P(E_i^- | E_i^+)} \\ &= \frac{P(E_i^- | X) P(X | E_i^+)}{P(E_i^- | E_i^+)} \end{aligned}$$



Decomposing the probabilities, cont'd

$$\begin{aligned}
 P(X_i | E) &= P(X_i | E_i^-, E_i^+) \\
 &= \frac{P(E_i^- | X, E_i^+) P(X | E_i^+)}{P(E_i^- | E_i^+)} \\
 &= \frac{P(E_i^- | X) P(X | E_i^+)}{P(E_i^- | E_i^+)} \\
 &= \alpha \pi(X_i) \lambda(X_i)
 \end{aligned}$$



These are each vectors.
Size of vector = # values of X_i

Where:

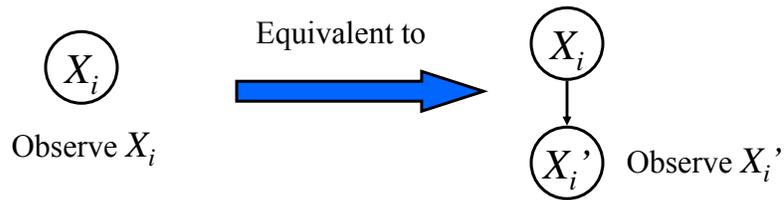
- α is a constant independent of X_i
- $\pi(X_i) = P(X_i | E_i^+)$
- $\lambda(X_i) = P(E_i^- | X_i)$

Using the decomposition for inference

- We can use this decomposition to do inference as follows. First, compute $\lambda(X_i) = P(E_i^- | X_i)$ for all X_i recursively, using the leaves of the tree as the base case.
- If X_i is a leaf:
 - If X_i is in E : $\lambda(X_i) = 1$ if X_i matches E , 0 otherwise
 - If X_i is not in E : E_i^- is the null set, so $P(E_i^- | X_i) = 1$ (constant)

Quick aside: “Virtual evidence”

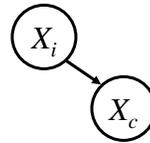
- For theoretical simplicity, but without loss of generality, let's assume that *all* variables in E (the evidence set) are leaves in the tree.
- Why can we do this WLOG:



Where $P(X_i' / X_i) = 1$ if $X_i' = X_i$, 0 otherwise

Calculating $\lambda(X_i)$ for non-leaves

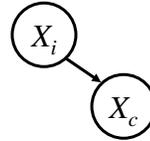
- Suppose X_i has one child, X_c .
- Then:



$$\lambda(X_i) = P(E_i^- | X_i) =$$

Calculating $\lambda(X_i)$ for non-leaves

- Suppose X_i has one child, X_c .

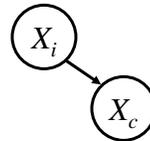


- Then:

$$\lambda(X_i) = P(E_i^- | X_i) = \sum_j P(E_i^-, X_c = j | X_i)$$

Calculating $\lambda(X_i)$ for non-leaves

- Suppose X_i has one child, X_c .

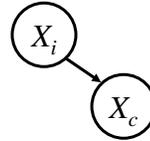


- Then:

$$\begin{aligned} \lambda(X_i) &= P(E_i^- | X_i) = \sum_j P(E_i^-, X_c = j | X_i) \\ &= \sum_j P(X_c = j | X_i) P(E_i^- | X_i, X_c = j) \end{aligned}$$

Calculating $\lambda(X_i)$ for non-leaves

- Suppose X_i has one child, X_c .



- Then:

$$\begin{aligned}
 \lambda(X_i) &= P(E_i^- | X_i) = \sum_j P(E_i^-, X_c = j | X_i) \\
 &= \sum_j P(X_c = j | X_i) P(E_i^- | X_i, X_c = j) \\
 &= \sum_j P(X_c = j | X_i) P(E_i^- | X_c = j) \\
 &= \sum_j P(X_c = j | X_i) \lambda(X_c = j)
 \end{aligned}$$

Calculating $\lambda(X_i)$ for non-leaves

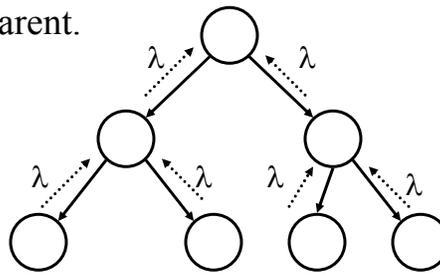
- Now, suppose X_i has a set of children, C .
- Since X_i *d-separates* each of its subtrees, the contribution of each subtree to $\lambda(X_i)$ is independent:

$$\begin{aligned}
 \lambda(X_i) &= P(E_i^- | X_i) = \prod_{X_j \in C} \lambda_j(X_i) \\
 &= \prod_{X_j \in C} \left[\sum_{X_j} P(X_j | X_i) \lambda(X_j) \right]
 \end{aligned}$$

where $\lambda_j(X_i)$ is the contribution to $P(E_i^- | X_i)$ of the part of the evidence lying in the subtree rooted at one of X_i 's children X_j .

We are now λ -happy

- So now we have a way to recursively compute all the $\lambda(X_i)$'s, starting from the root and using the leaves as the base case.
- If we want, we can think of each node in the network as an autonomous processor that passes a little “ λ message” to its parent.



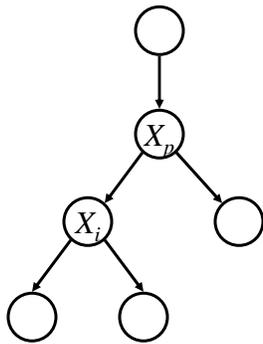
The other half of the problem

- Remember, $P(X_i|E) = \alpha\pi(X_i)\lambda(X_i)$. Now that we have all the $\lambda(X_i)$'s, what about the $\pi(X_i)$'s?

$$\pi(X_i) = P(X_i | E_i^+).$$

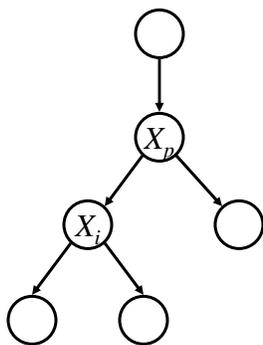
- What about the root of the tree, X_r ? In that case, E_r^+ is the null set, so $\pi(X_r) = P(X_r)$. No sweat. Since we also know $\lambda(X_r)$, we can compute the final $P(X_r)$.
- So for an arbitrary X_i with parent X_p , let's inductively assume we know $\pi(X_p)$ and/or $P(X_p/E)$. How do we get $\pi(X_i)$?

Computing $\pi(X_i)$



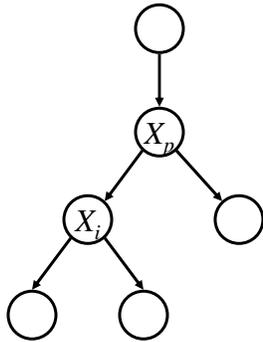
$$\pi(X_i) = P(X_i | E_i^+) =$$

Computing $\pi(X_i)$



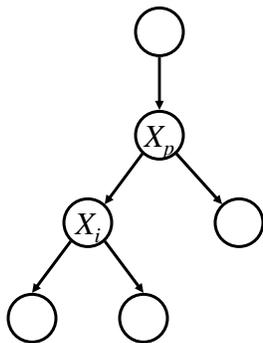
$$\pi(X_i) = P(X_i | E_i^+) = \sum_j P(X_i, X_p = j | E_i^+)$$

Computing $\pi(X_i)$



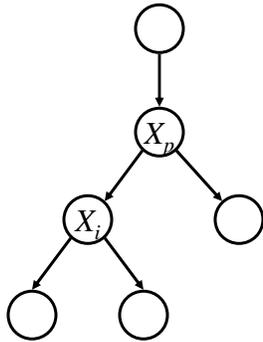
$$\begin{aligned}\pi(X_i) &= P(X_i | E_i^+) = \sum_j P(X_i, X_p = j | E_i^+) \\ &= \sum_j P(X_i | X_p = j, E_i^+) P(X_p = j | E_i^+)\end{aligned}$$

Computing $\pi(X_i)$



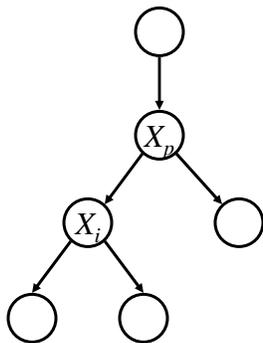
$$\begin{aligned}\pi(X_i) &= P(X_i | E_i^+) = \sum_j P(X_i, X_p = j | E_i^+) \\ &= \sum_j P(X_i | X_p = j, E_i^+) P(X_p = j | E_i^+) \\ &= \sum_j P(X_i | X_p = j) P(X_p = j | E_i^+)\end{aligned}$$

Computing $\pi(X_i)$



$$\begin{aligned}
 \pi(X_i) &= P(X_i | E_i^+) = \sum_j P(X_i, X_p = j | E_i^+) \\
 &= \sum_j P(X_i | X_p = j, E_i^+) P(X_p = j | E_i^+) \\
 &= \sum_j P(X_i | X_p = j) P(X_p = j | E_i^+) \\
 &= \sum_j P(X_i | X_p = j) \frac{P(X_p = j | E)}{\lambda_i(X_p = j)}
 \end{aligned}$$

Computing $\pi(X_i)$

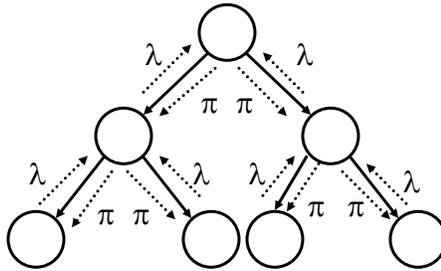


$$\begin{aligned}
 \pi(X_i) &= P(X_i | E_i^+) = \sum_j P(X_i, X_p = j | E_i^+) \\
 &= \sum_j P(X_i | X_p = j, E_i^+) P(X_p = j | E_i^+) \\
 &= \sum_j P(X_i | X_p = j) P(X_p = j | E_i^+) \\
 &= \sum_j P(X_i | X_p = j) \frac{P(X_p = j | E)}{\lambda_i(X_p = j)} \\
 &= \sum_j P(X_i | X_p = j) \pi_i(X_p = j)
 \end{aligned}$$

Where $\pi_i(X_p)$ is defined as $\frac{P(X_p | E)}{\lambda_i(X_p)}$

We're done. Yay!

- Thus we can compute all the $\pi(X_i)$'s, and, in turn, all the $P(X_i|E)$'s.
- Can think of nodes as autonomous processors passing λ and π messages to their neighbors

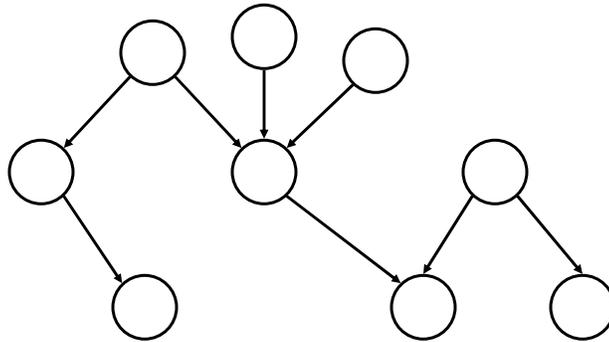


Conjunctive queries

- What if we want, e.g., $P(A, B | C)$ instead of just marginal distributions $P(A | C)$ and $P(B | C)$?
- Just use chain rule:
 - $P(A, B | C) = P(A | C) P(B | A, C)$
 - Each of the latter probabilities can be computed using the technique just discussed.

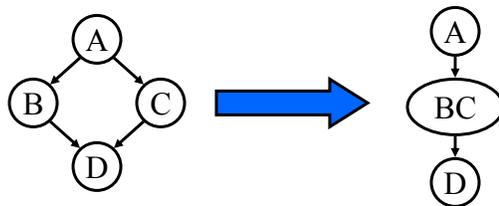
Polytrees

- Technique can be generalized to *polytrees*: undirected versions of the graphs are still trees, but nodes can have more than one parent

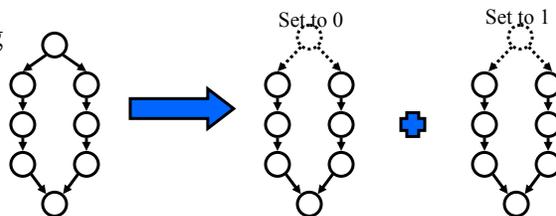


Dealing with cycles

- Can deal with undirected cycles in graph by
 - clustering variables together

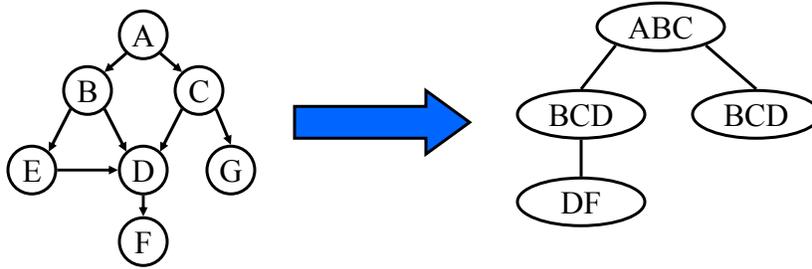


- Conditioning



Join trees

- Arbitrary Bayesian network can be transformed via some evil graph-theoretic magic into a *join tree* in which a similar method can be employed.



In the worst case the join tree nodes must take on exponentially many combinations of values, but often works well in practice