

Belief Propagation in Bayesian Networks

Céline Comte

NOKIA Bell Labs



Reading Group “Network Theory”
November 5, 2018

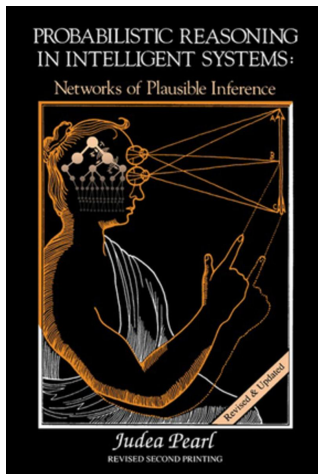
Introduction

(First-order) logic

Represent causal relations between variables by a directed acyclic graph

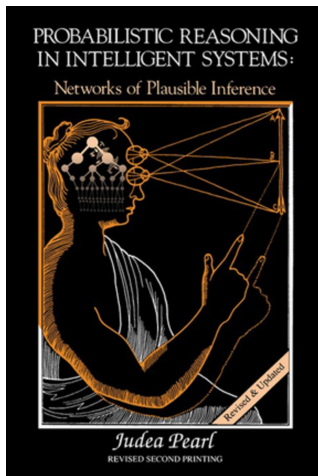
Probabilities

Weight these causal relations by probabilities that implicitly account for non-represented variables



Introduction

“Belief networks are directed acyclic graphs in which the nodes represent propositions (or variables), the arcs signify direct dependencies between the linked propositions, and the strengths of these dependencies are quantified by conditional probabilities” (Pearl, 1986)



Bayesian networks are also ...

- A memory-efficient way of storing a PMF
- Based on simple probability rules
(more details in a few slides)
- Inspired by human causal reasoning (Pearl, 1986, 1988)
- Used for decision taking if a utility function is provided
- Applied in many fields: medicine diagnoses, turbo-codes, (programming) language detection, ...
- Related to other models: Markov random fields, Markov chains, hidden Markov models, ...

References: Pearl's articles and book

- [J. Pearl \(1982\)](#). “Reverend Bayes on Inference Engines: A Distributed Hierarchical Approach”. In: [AAAI'82](#)
 - Belief propagation in causal trees
- [J. Pearl \(1986\)](#). “Fusion, propagation, and structuring in belief networks”. In: [Artificial Intelligence](#)
 - Belief propagation in causal trees and polytrees
- [J. Pearl \(1988\)](#). Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. [Morgan Kaufmann](#)
 - A complete reference
(thanks Achille for providing me with this book)

References: Textbooks

- T. D. Nielsen and F. V. Jensen (2007). Bayesian Networks and Decision Graphs. Springer-Verlag
→ A lot of examples in Chapters 2 and 3
- D. Koller, N. Friedman, and F. Bach (2009). Probabilistic Graphical Models: Principles and Techniques. MIT Press
- M. Jordan (Last modified in 2015). An Introduction to Probabilistic Graphical Models.
→ Definition and belief propagation (thanks Nathan for pointing this reference)

Outline

Reminders on probability theory

Bayesian networks

Belief propagation in trees

Belief propagation in polytrees

Outline

Reminders on probability theory

Bayesian networks

Belief propagation in trees

Belief propagation in polytrees

Independence and conditional independence

Remark: We work exclusively with discrete random variables

Independence and conditional independence

Remark: We work exclusively with discrete random variables

- A and B are **marginally independent** (written $A \perp B$) if one of these three equivalent conditions is satisfied:
 - $P(A, B) = P(A)P(B)$
 - $P(A | B) = P(A)$
 - $P(B | A) = P(B)$

Independence and conditional independence

Remark: We work exclusively with discrete random variables

- A and B are **marginally independent** (written $A \perp B$) if one of these three equivalent conditions is satisfied:
 - $P(A, B) = P(A)P(B)$
 - $P(A | B) = P(A)$
 - $P(B | A) = P(B)$
- A and B are **conditionally independent** given C (written $A \perp B | C$) if one of these three equivalent conditions is satisfied:
 - $P(A, B | C) = P(A | C)P(B | C)$
 - $P(A | B, C) = P(A | C)$
 - $P(B | A, C) = P(B | C)$

Useful rules

- **The chain rule of probabilities**

If A_1, \dots, A_n are random variables, we have

$$P(A_1, \dots, A_n) = P(A_1) \times P(A_2 | A_1) \times P(A_3 | A_1, A_2) \\ \times \dots \times P(A_n | A_1, \dots, A_{n-1})$$

Useful rules

- **The chain rule of probabilities**

If A_1, \dots, A_n are random variables, we have

$$P(A_1, \dots, A_n) = P(A_1) \times P(A_2 | A_1) \times P(A_3 | A_1, A_2) \\ \times \dots \times P(A_n | A_1, \dots, A_{n-1})$$

- **Law of total probability**

If A and B are two random variables,

$$P(B) = \sum_A P(B | A)P(A)$$

Useful rules

- **Bayes' rule**

If A and B are two random variables,

$$P(B | A) = \frac{P(A | B)P(B)}{P(A)}$$

We can see $P(A)$ as a **normalizing constant**: we can first compute $P(B | A) \propto P(A | B)P(B)$ for each value of B and then normalize to obtain $P(B | A)$ without computing $P(A)$

Glossary

- **Belief** in a random variable (**conviction** in french)
Marginal distribution of this random variable
(given the value of some observed variables)

Glossary

- **Belief** in a random variable (**conviction** in french)
Marginal distribution of this random variable
(given the value of some observed variables)
- **Observe** a random variable

Glossary

- **Belief** in a random variable (**conviction** in french)
Marginal distribution of this random variable
(given the value of some observed variables)
- **Observe** a random variable
- **Evidence** (piece of evidence)
The set of random variables that have been observed

Outline

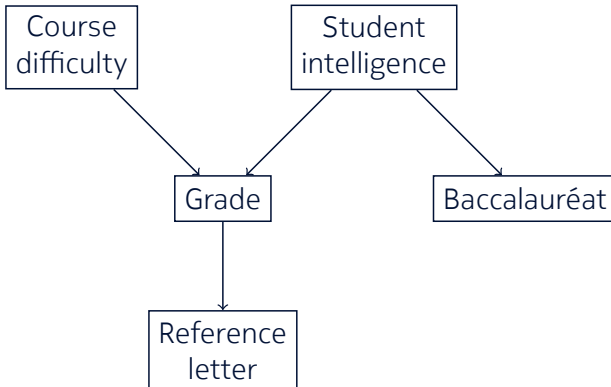
Reminders on probability theory

Bayesian networks

Belief propagation in trees

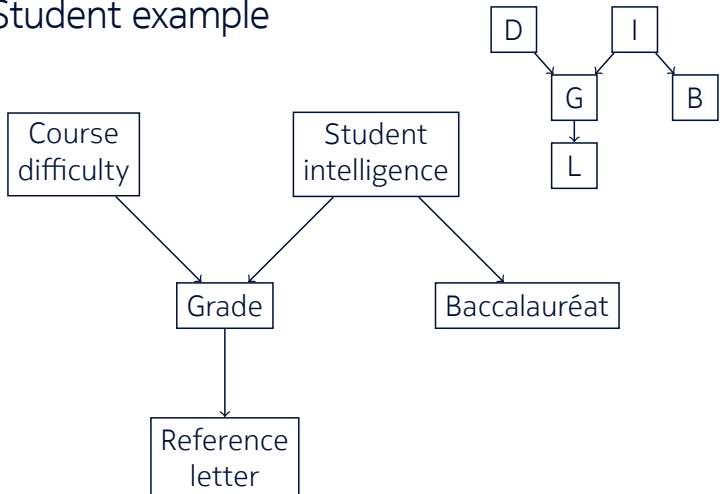
Belief propagation in polytrees

The Student example



Borrowed from (Koller, Friedman, and Bach, 2009)

The Student example

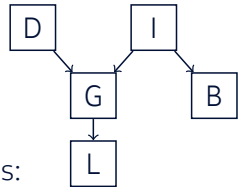


Borrowed from (Koller, Friedman, and Bach, 2009)

The Student example

- **Local Markov property**

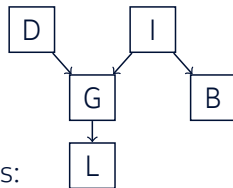
Each node is conditionally independent of its non-descendants given its parents:



$$D \perp\!\!\!\perp \{I, B\}, \quad I \perp\!\!\!\perp D, \quad G \perp\!\!\!\perp B \mid \{D, I\},$$

$$B \perp\!\!\!\perp \{D, G, L\} \mid I, \quad L \perp\!\!\!\perp \{D, I, B\} \mid G$$

The Student example



- **Local Markov property**

Each node is conditionally independent of its non-descendants given its parents:

$$D \perp\!\!\!\perp \{I, B\}, \quad I \perp\!\!\!\perp D, \quad G \perp\!\!\!\perp B \mid \{D, I\},$$

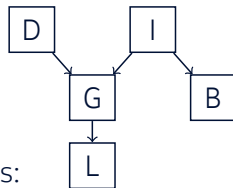
$$B \perp\!\!\!\perp \{D, G, L\} \mid I, \quad L \perp\!\!\!\perp \{D, I, B\} \mid G$$

- **Chain rule of Bayesian networks**

By the chain rule of probabilities:

$$\begin{aligned} P(D, I, G, B, L) &= P(D)P(I \mid D)P(G \mid D, I)P(B \mid D, I, G)P(L \mid D, I, G, B), \\ &= P(D)P(I)P(G \mid D, I)P(B \mid I)P(L \mid G) \end{aligned}$$

The Student example



- **Local Markov property**

Each node is conditionally independent of its non-descendants given its parents:

$$D \perp\!\!\!\perp \{I, B\}, \quad I \perp\!\!\!\perp D, \quad G \perp\!\!\!\perp B \mid \{D, I\},$$

$$B \perp\!\!\!\perp \{D, G, L\} \mid I, \quad L \perp\!\!\!\perp \{D, I, B\} \mid G$$

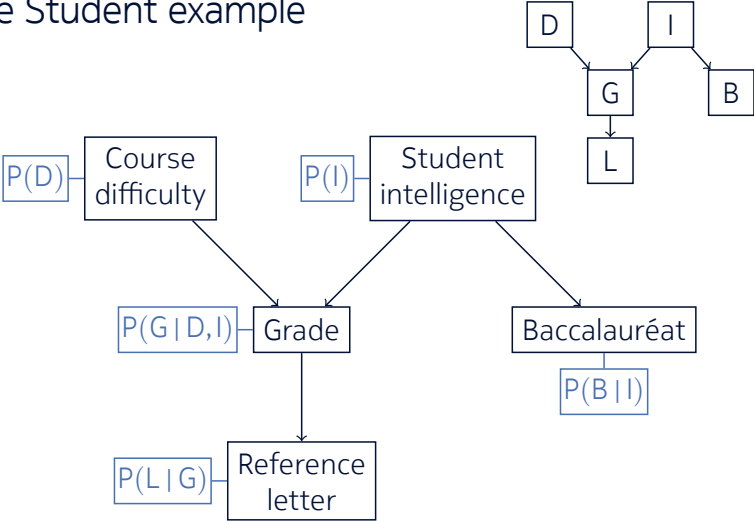
- **Chain rule of Bayesian networks**

By the chain rule of probabilities:

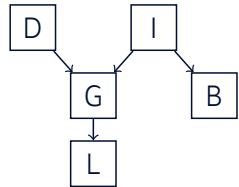
$$\begin{aligned} P(D, I, G, B, L) &= P(D)P(I \mid D)P(G \mid D, I)P(B \mid D, I, G)P(L \mid D, I, G, B), \\ &= P(D)P(I)P(G \mid D, I)P(B \mid I)P(L \mid G) \end{aligned}$$

These two definitions are equivalent

The Student example



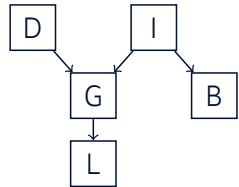
Bayesian networks in general



Described by

- A **directed acyclic graph**
 - Nodes \sim (discrete) random variables X_1, \dots, X_n
 - Arrows \sim conditional (in)dependencies
- Local **conditional probability tables (CPT)**
 - $P(X_i | \text{parents}(X_i))$ for each node X_i

Bayesian networks in general



Two equivalent definitions

- **Local Markov property**

Each node is conditionally independent of its non-descendants given its parents

- **Chain rule of Bayesian networks**

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

Proof of the equivalence: Corollary 4 p.20 of (Pearl, 1988)

Base case: serial connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X)P(Y \mid X)P(Z \mid Y)$$

Base case: serial connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X)P(Y \mid X)P(Z \mid Y)$$

- Interpretation: chain of causality
X “causes” Y that “causes” Z

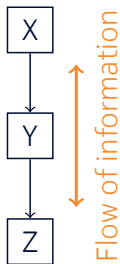
Base case: serial connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X)P(Y|X)P(Z|Y)$$

- Interpretation: chain of causality
X “causes” Y that “causes” Z
- Information can flow between X and Z through Y (that is, observing X changes our belief in Z and vice versa), unless Y is observed

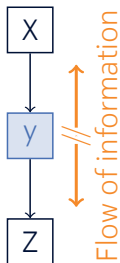
Base case: serial connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X)P(Y|X)P(Z|Y)$$

- Interpretation: chain of causality
X “causes” Y that “causes” Z
- Information can flow between X and Z through Y (that is, observing X changes our belief in Z and vice versa), unless Y is observed

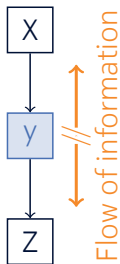
Base case: serial connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X)P(Y|X)P(Z|Y)$$

- Interpretation: chain of causality
X “causes” Y that “causes” Z
- Information can flow between X and Z through Y (that is, observing X changes our belief in Z and vice versa), unless Y is observed

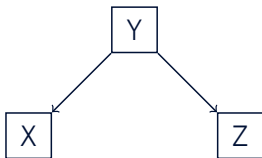
Base case: serial connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X)P(Y \mid X)P(Z \mid Y)$$

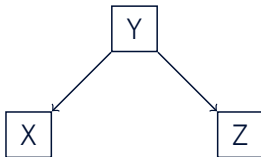
- Interpretation: chain of causality
X “causes” Y that “causes” Z
- Information can flow between X and Z through Y (that is, observing X changes our belief in Z and vice versa), unless Y is observed
- Example: Markov chains

Base case: diverging connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X \mid Y)P(Y)P(Z \mid Y)$$

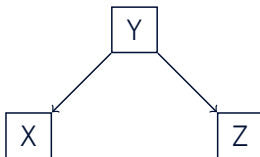
Base case: diverging connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X \mid Y)P(Y)P(Z \mid Y)$$

- Interpretation: a single root cause Y with two observable consequences X and Z

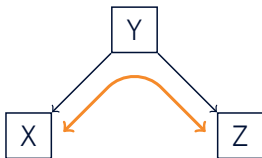
Base case: diverging connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X \mid Y)P(Y)P(Z \mid Y)$$

- Interpretation: a single root cause Y with two observable consequences X and Z
- Information can flow between X and Z through Y, unless Y is observed

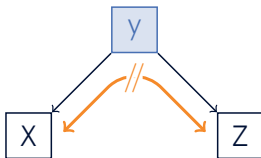
Base case: diverging connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X \mid Y)P(Y)P(Z \mid Y)$$

- Interpretation: a single root cause Y with two observable consequences X and Z
- Information can flow between X and Z through Y , unless Y is observed

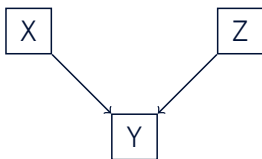
Base case: diverging connection



$$X \perp\!\!\!\perp Z \mid Y \quad P(X, Y, Z) = P(X \mid Y)P(Y)P(Z \mid Y)$$

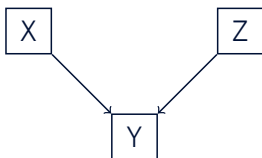
- Interpretation: a single root cause Y with two observable consequences X and Z
- Information can flow between X and Z through Y , unless Y is observed

Base case: converging connection



$$X \perp\!\!\!\perp Z \quad P(X, Y, Z) = P(X)P(Y | X, Z)P(Z)$$

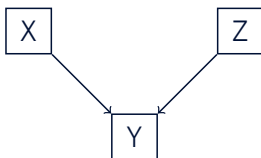
Base case: converging connection



$$X \perp\!\!\!\perp Z \quad P(X, Y, Z) = P(X)P(Y | X, Z)P(Z)$$

- Interpretation: two possible explanations X and Z for an observed consequence Y

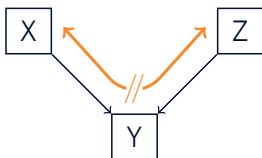
Base case: converging connection



$$X \perp\!\!\!\perp Z \quad P(X, Y, Z) = P(X)P(Y | X, Z)P(Z)$$

- Interpretation: two possible explanations X and Z for an observed consequence Y
- “Explaining away” effect: information cannot flow between X and Z, unless Y is observed

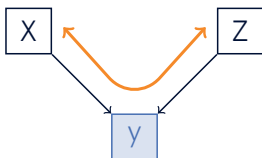
Base case: converging connection



$$X \perp\!\!\!\perp Z \quad P(X, Y, Z) = P(X)P(Y | X, Z)P(Z)$$

- Interpretation: two possible explanations X and Z for an observed consequence Y
- “Explaining away” effect: information cannot flow between X and Z, unless Y is observed

Base case: converging connection

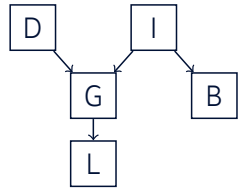


$$X \perp\!\!\!\perp Z \quad P(X, Y, Z) = P(X)P(Y | X, Z)P(Z)$$

- Interpretation: two possible explanations X and Z for an observed consequence Y
- “Explaining away” effect: information cannot flow between X and Z, unless Y is observed

Implied independencies

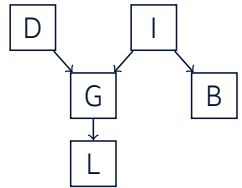
Similar to the “Strong Markov property”
of Markov chains



Implied independencies

Similar to the “Strong Markov property”
of Markov chains

Which are correct?

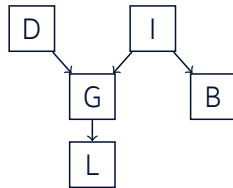


Implied independencies

Similar to the “Strong Markov property” of Markov chains

Which are correct?

- ① $G \perp\!\!\!\perp B$?

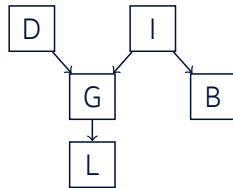


Implied independencies

Similar to the “Strong Markov property”
of Markov chains

Which are correct?

- ① $G \perp\!\!\!\perp B$? No

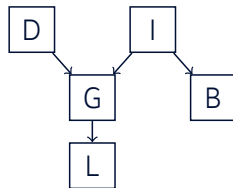


Implied independencies

Similar to the “Strong Markov property”
of Markov chains

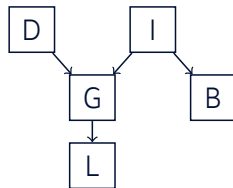
Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$?



Implied independencies

Similar to the “Strong Markov property” of Markov chains

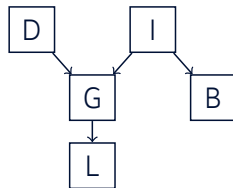


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No

Implied independencies

Similar to the “Strong Markov property” of Markov chains

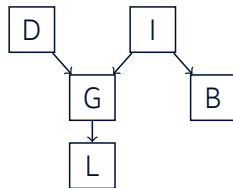


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$?

Implied independencies

Similar to the “Strong Markov property” of Markov chains

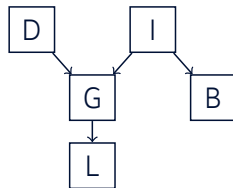


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$? No

Implied independencies

Similar to the “Strong Markov property” of Markov chains

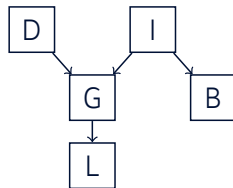


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$? No
- ④ $D \perp\!\!\!\perp B$?

Implied independencies

Similar to the “Strong Markov property” of Markov chains

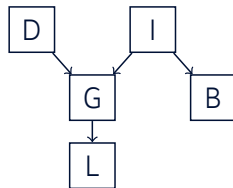


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$? No
- ④ $D \perp\!\!\!\perp B$? Yes (by the local Markov property applied to D)

Implied independencies

Similar to the “Strong Markov property” of Markov chains

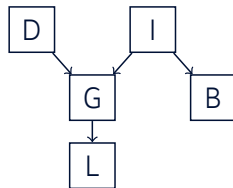


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$? No
- ④ $D \perp\!\!\!\perp B$? Yes (by the local Markov property applied to D)
- ⑤ $D \perp\!\!\!\perp B \mid G$?

Implied independencies

Similar to the “Strong Markov property” of Markov chains

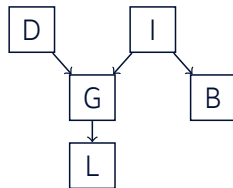


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$? No
- ④ $D \perp\!\!\!\perp B$? Yes (by the local Markov property applied to D)
- ⑤ $D \perp\!\!\!\perp B \mid G$? No (“explaining away” effect)

Implied independencies

Similar to the “Strong Markov property” of Markov chains

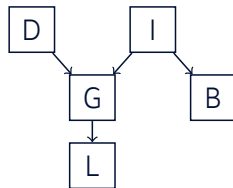


Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$? No
- ④ $D \perp\!\!\!\perp B$? Yes (by the local Markov property applied to D)
- ⑤ $D \perp\!\!\!\perp B \mid G$? No (“explaining away” effect)
- ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$?

Implied independencies

Similar to the “Strong Markov property” of Markov chains



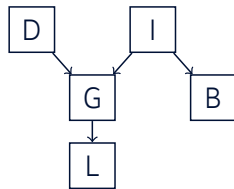
Which are correct?

- ① $G \perp\!\!\!\perp B$? No
- ② $B \perp\!\!\!\perp L$? No
- ③ $D \perp\!\!\!\perp L$? No
- ④ $D \perp\!\!\!\perp B$? Yes (by the local Markov property applied to D)
- ⑤ $D \perp\!\!\!\perp B \mid G$? No (“explaining away” effect)
- ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$? Yes

Implied independencies

Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

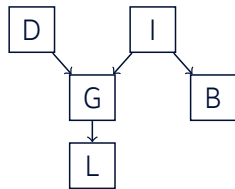
$P(D, B \mid I, G)$



Implied independencies

Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

$$P(D, B \mid I, G)$$



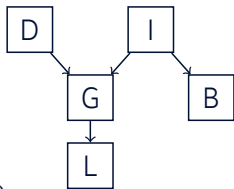
$$= P(D \mid G, I)P(B \mid I, G)$$

□

Implied independencies

Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

$$P(D, B \mid I, G) = \frac{P(G \mid D, I, B)P(D, B \mid I)}{P(G \mid I)} \quad (\text{Bayes' rule})$$



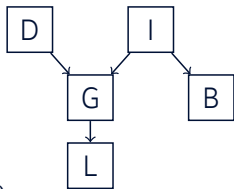
$$= P(D \mid G, I)P(B \mid I, G)$$

□

Implied independencies

Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

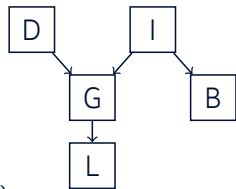
$$\begin{aligned} P(D, B \mid I, G) &= \frac{P(G \mid D, I, B)P(D, B \mid I)}{P(G \mid I)} && \text{(Bayes' rule)} \\ &= \frac{P(G \mid D, I)P(D, B \mid I)}{P(G \mid I)} && \text{(local Markov property applied to G)} \end{aligned}$$



$$= P(D \mid G, I)P(B \mid I, G)$$

□

Implied independencies



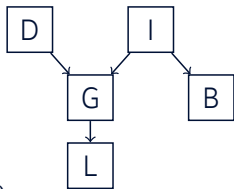
Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

$$\begin{aligned} P(D, B \mid I, G) &= \frac{P(G \mid D, I, B)P(D, B \mid I)}{P(G \mid I)} && \text{(Bayes' rule)} \\ &= \frac{P(G \mid D, I)P(D, B \mid I)}{P(G \mid I)} && \text{(local Markov property applied to G)} \\ &= \frac{P(G \mid D, I)P(D \mid I)P(B \mid D, I)}{P(G \mid I)} && \text{(definition of conditional probabilities)} \end{aligned}$$

$$= P(D \mid G, I)P(B \mid I, G)$$

□

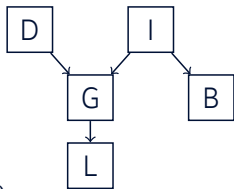
Implied independencies



Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

$$\begin{aligned} P(D, B \mid I, G) &= \frac{P(G \mid D, I, B)P(D, B \mid I)}{P(G \mid I)} && \text{(Bayes' rule)} \\ &= \frac{P(G \mid D, I)P(D, B \mid I)}{P(G \mid I)} && \text{(local Markov property applied to G)} \\ &= \frac{P(G \mid D, I)P(D \mid I)P(B \mid D, I)}{P(G \mid I)} && \text{(definition of conditional probabilities)} \\ &= \frac{P(G \mid D, I)P(D \mid I)}{P(G \mid I)} P(B \mid D, I) \\ &= P(D \mid G, I)P(B \mid I, G) \quad \square \end{aligned}$$

Implied independencies

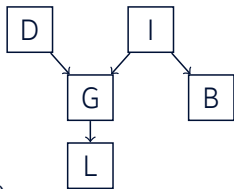


Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

$$\begin{aligned} P(D, B \mid I, G) &= \frac{P(G \mid D, I, B)P(D, B \mid I)}{P(G \mid I)} && \text{(Bayes' rule)} \\ &= \frac{P(G \mid D, I)P(D, B \mid I)}{P(G \mid I)} && \text{(local Markov property applied to G)} \\ &= \frac{P(G \mid D, I)P(D \mid I)P(B \mid D, I)}{P(G \mid I)} && \text{(definition of conditional probabilities)} \\ &= \frac{P(G \mid D, I)P(D \mid I)}{P(G \mid I)} P(B \mid D, I) \\ &= P(D \mid G, I)P(B \mid D, I) && \text{(Bayes' rule)} \\ &= P(D \mid G, I)P(B \mid I, G) \end{aligned}$$

□

Implied independencies



Proof of ⑥ $D \perp\!\!\!\perp B \mid \{I, G\}$

$$\begin{aligned} P(D, B \mid I, G) &= \frac{P(G \mid D, I, B)P(D, B \mid I)}{P(G \mid I)} && \text{(Bayes' rule)} \\ &= \frac{P(G \mid D, I)P(D, B \mid I)}{P(G \mid I)} && \text{(local Markov property applied to G)} \\ &= \frac{P(G \mid D, I)P(D \mid I)P(B \mid D, I)}{P(G \mid I)} && \text{(definition of conditional probabilities)} \\ &= \frac{P(G \mid D, I)P(D \mid I)}{P(G \mid I)}P(B \mid D, I) \\ &= P(D \mid G, I)P(B \mid D, I) && \text{(Bayes' rule)} \\ &= P(D \mid G, I)P(B \mid I, G) && \text{(local Markov property applied to B)} \quad \square \end{aligned}$$

Memory and time complexity

Parameters

Memory and time complexity

Parameters

- n number of random variables
(typically, $n \sim$ hundreds or thousands)

Memory and time complexity

Parameters

- n number of random variables
(typically, $n \sim$ hundreds or thousands)
- r number of values each variable can take

Memory and time complexity

Parameters

n number of random variables

(typically, $n \sim$ hundreds or thousands)

r number of values each variable can take

d^\dagger maximum number of parents of a node

Memory and time complexity

Parameters

- n number of random variables
(typically, $n \sim$ hundreds or thousands)
- r number of values each variable can take
- d^\dagger maximum number of parents of a node

Memory complexity

Memory and time complexity

Parameters

- n number of random variables
(typically, $n \sim$ hundreds or thousands)
- r number of values each variable can take
- d^\dagger maximum number of parents of a node

Memory complexity

- If we store the probability distribution: $O(r^n)$ entries

Memory and time complexity

Parameters

- n number of random variables
(typically, $n \sim$ hundreds or thousands)
- r number of values each variable can take
- d^{\dagger} maximum number of parents of a node

Memory complexity

- If we store the probability distribution: $O(r^n)$ entries
- If we store the node parents and the conditional probability tables: $O(n(r + r^{d^{\dagger}})) = O(nr^{d^{\dagger}})$ entries

Memory and time complexity

Parameters

- n number of random variables
(typically, $n \sim$ hundreds or thousands)
- r number of values each variable can take
- d^\dagger maximum number of parents of a node

Memory complexity

- If we store the probability distribution: $O(r^n)$ entries
- If we store the node parents and the conditional probability tables: $O(n(r + r^{d^\dagger})) = O(nr^{d^\dagger})$ entries

What about the time complexity?

Inference

“A guess that you make or an opinion that you form based on the information that you have” (Cambridge dictionary)

Inference

“A guess that you make or an opinion that you form based on the information that you have” (Cambridge dictionary)

→ Bayesian networks: compute or update the belief in each variable given some evidence

Inference

“A guess that you make or an opinion that you form based on the information that you have” (Cambridge dictionary)

→ Bayesian networks: compute or update the belief in each variable given some evidence

Belief propagation, a.k.a. **sum-product message passing**:

Propagate the information through the network, starting from the evidence node(s)

Inference

“A guess that you make or an opinion that you form based on the information that you have” (Cambridge dictionary)

→ Bayesian networks: compute or update the belief in each variable given some evidence

Belief propagation, a.k.a. **sum-product message passing**:

Propagate the information through the network, starting from the evidence node(s)

- Each variable is a “separate processor” (a neuron?) that knows its own CPT and the messages received from its direct neighbors (Pearl, 1982)

Inference

“A guess that you make or an opinion that you form based on the information that you have” (Cambridge dictionary)

→ Bayesian networks: compute or update the belief in each variable given some evidence

Belief propagation, a.k.a. **sum-product message passing**:

Propagate the information through the network, starting from the evidence node(s)

- Each variable is a “separate processor” (a neuron?) that knows its own CPT and the messages received from its direct neighbors (Pearl, 1982)
- Dynamic programming

Outline

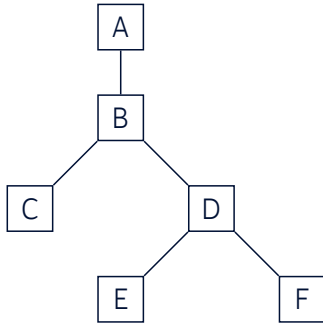
Reminders on probability theory

Bayesian networks

Belief propagation in trees

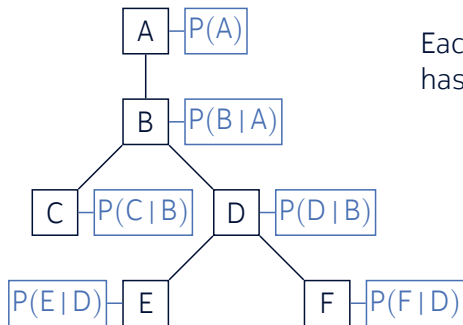
Belief propagation in polytrees

Tree Bayesian network



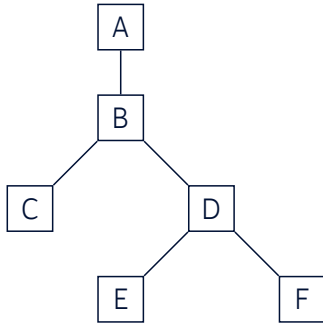
Each node (except the root) has at most one parent

Tree Bayesian network



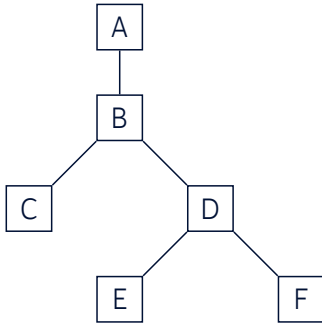
Each node (except the root) has at most one parent

Tree Bayesian network



Each node (except the root) has at most one parent

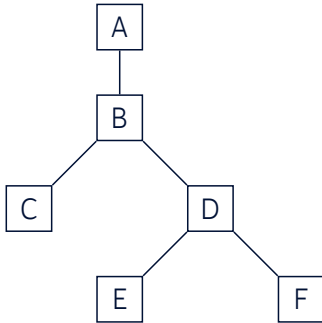
Tree Bayesian network



Each node (except the root) has at most one parent

Each node **separates** the tree: its non-descendants and the subtrees rooted at each of its children are conditionally independent given this node

Tree Bayesian network

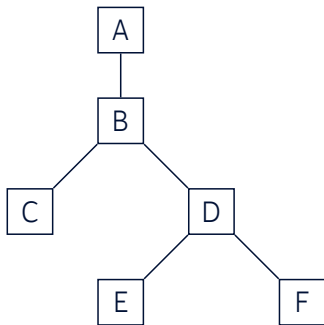


Each node (except the root) has at most one parent

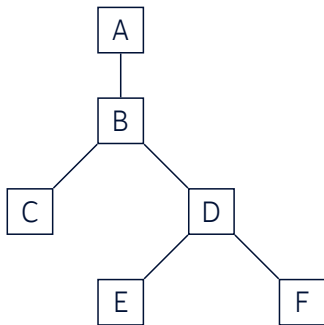
Each node **separates** the tree: its non-descendants and the subtrees rooted at each of its children are conditionally independent given this node

Remark: We will explain the propagation algorithm on this toy example borrowed from (Pearl, 1988)

No evidence

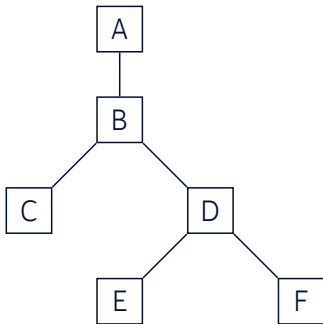


No evidence



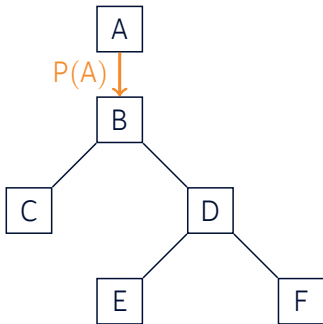
- $P(A)$: parameter

No evidence



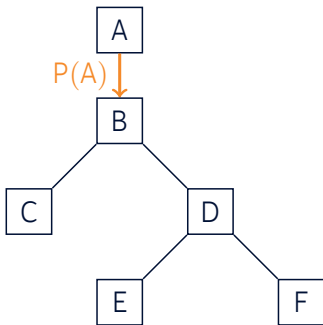
- $P(A)$: parameter
- $P(B) = \sum_A P(B|A)P(A)$

No evidence



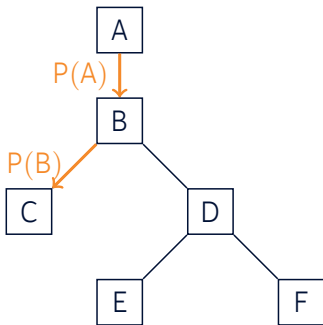
- $P(A)$: parameter
- $P(B) = \sum_A P(B|A)P(A)$

No evidence



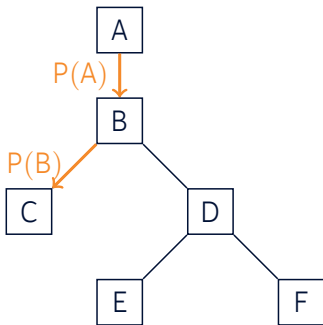
- $P(A)$: parameter
- $P(B) = \sum_A P(B|A)P(A)$
- $P(C) = \sum_B P(C|B)P(B)$

No evidence



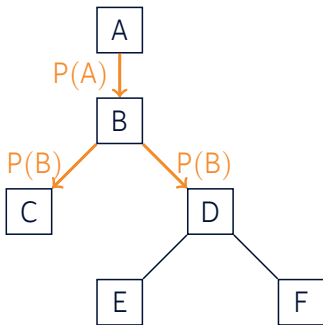
- $P(A)$: parameter
- $P(B) = \sum_A P(B|A)P(A)$
- $P(C) = \sum_B P(C|B)P(B)$

No evidence



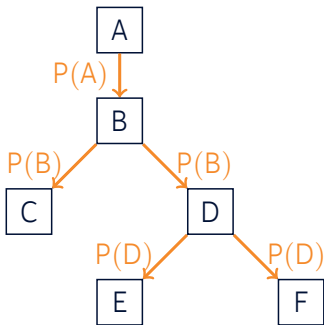
- $P(A)$: parameter
- $P(B) = \sum_A P(B|A)P(A)$
- $P(C) = \sum_B P(C|B)P(B)$
- $P(D) = \sum_B P(D|B)P(B)$

No evidence



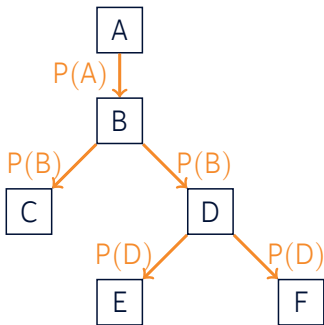
- $P(A)$: parameter
- $P(B) = \sum_A P(B|A)P(A)$
- $P(C) = \sum_B P(C|B)P(B)$
- $P(D) = \sum_B P(D|B)P(B)$

No evidence



- $P(A)$: parameter
- $P(B) = \sum_A P(B|A)P(A)$
- $P(C) = \sum_B P(C|B)P(B)$
- $P(D) = \sum_B P(D|B)P(B)$

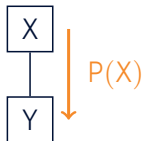
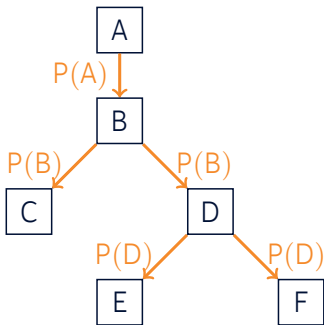
No evidence



- $P(A)$: parameter
- $P(B) = \sum_A P(B | A)P(A)$
- $P(C) = \sum_B P(C | B)P(B)$
- $P(D) = \sum_B P(D | B)P(B)$

Top-down propagation
Complexity $O(nr^2)$

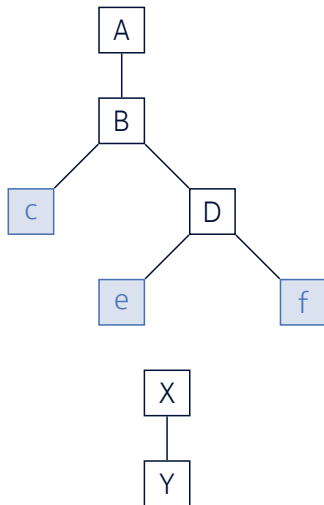
No evidence



- $P(A)$: parameter
- $P(B) = \sum_A P(B | A)P(A)$
- $P(C) = \sum_B P(C | B)P(B)$
- $P(D) = \sum_B P(D | B)P(B)$

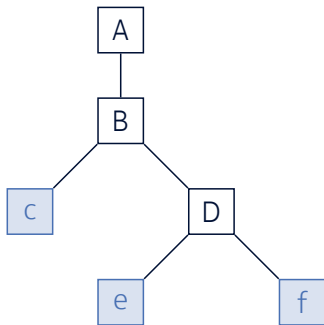
Top-down propagation
Complexity $O(nr^2)$

Three pieces of evidence



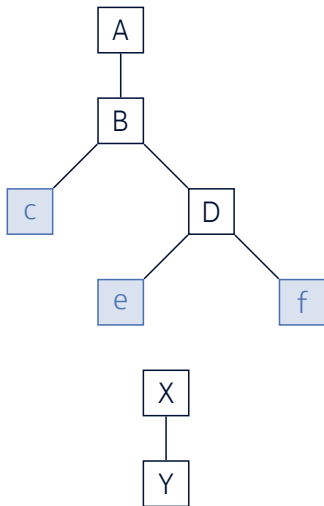
- **Evidence:** We observe that $C = c$, $E = e$, and $F = f$

Three pieces of evidence



- **Evidence:** We observe that $C = c$, $E = e$, and $F = f$
- **Objective:** Compute the belief $BEL(X) = P(X | c, e, f)$ of each node X

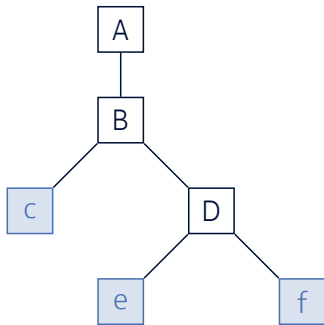
Three pieces of evidence



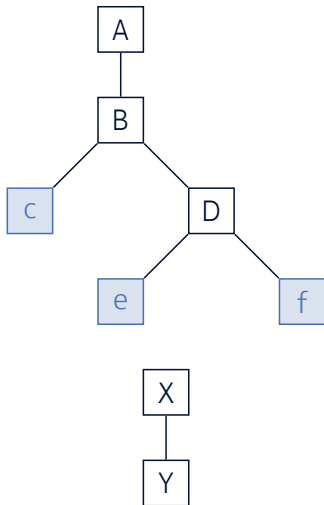
- **Evidence:** We observe that $C = c$, $E = e$, and $F = f$
- **Objective:** Compute the belief $BEL(X) = P(X | c, e, f)$ of each node X
- **Principle:** Propagate the information through the network, starting from the evidence nodes

Causal and diagnostic support

By Bayes' rule:



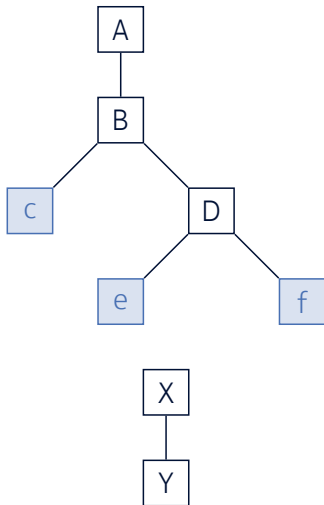
Causal and diagnostic support



By Bayes' rule:

- $$P(D | c, e, f) = \frac{P(e, f | D, c)P(D | c)}{P(e, f | c)}$$

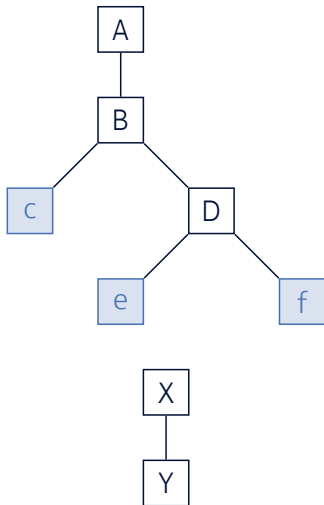
Causal and diagnostic support



By Bayes' rule:

- $$P(D | c, e, f) = \frac{P(e, f | D, c) P(D | c)}{P(e, f | c)}$$
$$= \frac{P(e, f | D) P(D | c)}{P(e, f | c)}$$

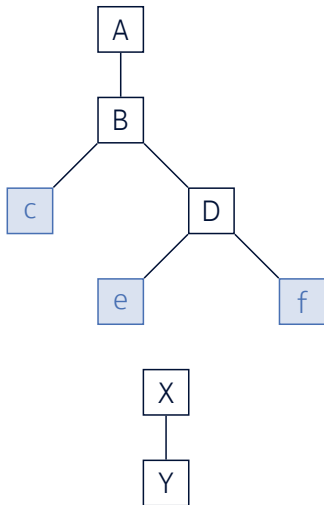
Causal and diagnostic support



By Bayes' rule:

- $$P(D | c, e, f) = \frac{P(e, f | D, c) P(D | c)}{P(e, f | c)}$$
$$= \frac{P(e, f | D) P(D | c)}{P(e, f | c)}$$
$$\propto P(e, f | D) P(D | c)$$

Causal and diagnostic support

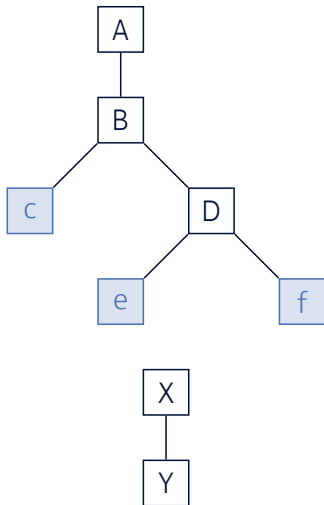


By Bayes' rule:

$$\bullet P(B | c, e, f) = \frac{P(c, e, f | B)P(B)}{P(c, e, f)}$$

$$\begin{aligned}\bullet P(D | c, e, f) &= \frac{P(e, f | D, c)P(D | c)}{P(e, f | c)} \\ &= \frac{P(e, f | D)P(D | c)}{P(e, f | c)} \\ &\propto P(e, f | D)P(D | c)\end{aligned}$$

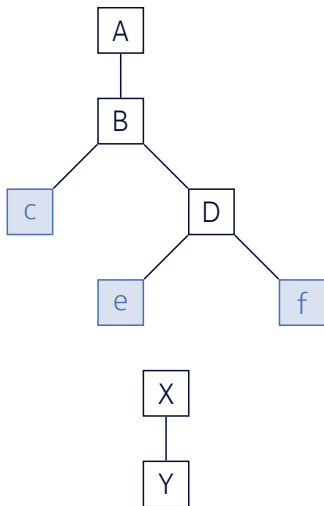
Causal and diagnostic support



By Bayes' rule:

- $$P(B | c, e, f) = \frac{P(c, e, f | B)P(B)}{P(c, e, f)} \propto P(c, e, f | B)P(B)$$
- $$P(D | c, e, f) = \frac{P(e, f | D, c)P(D | c)}{P(e, f | c)} = \frac{P(e, f | D)P(D | c)}{P(e, f | c)} \propto P(e, f | D)P(D | c)$$

Causal and diagnostic support



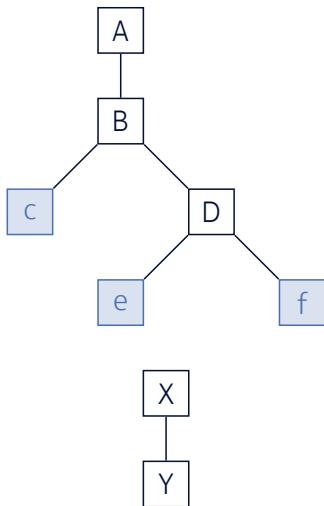
By Bayes' rule:

$$\bullet P(A | c, e, f) = \frac{P(c, e, f | A)P(A)}{P(c, e, f)}$$

$$\bullet P(B | c, e, f) = \frac{P(c, e, f | B)P(B)}{P(c, e, f)} \\ \propto P(c, e, f | B)P(B)$$

$$\bullet P(D | c, e, f) = \frac{P(e, f | D, c)P(D | c)}{P(e, f | c)} \\ = \frac{P(e, f | D)P(D | c)}{P(e, f | c)} \\ \propto P(e, f | D)P(D | c)$$

Causal and diagnostic support



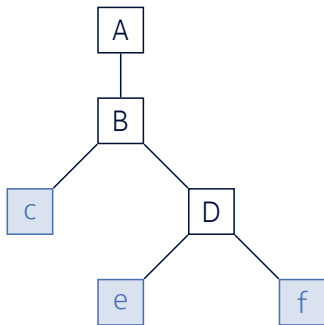
By Bayes' rule:

$$\begin{aligned} \bullet P(A | c, e, f) &= \frac{P(c, e, f | A)P(A)}{P(c, e, f)} \\ &\propto P(c, e, f | A)P(A) \end{aligned}$$

$$\begin{aligned} \bullet P(B | c, e, f) &= \frac{P(c, e, f | B)P(B)}{P(c, e, f)} \\ &\propto P(c, e, f | B)P(B) \end{aligned}$$

$$\begin{aligned} \bullet P(D | c, e, f) &= \frac{P(e, f | D, c)P(D | c)}{P(e, f | c)} \\ &= \frac{P(e, f | D)P(D | c)}{P(e, f | c)} \\ &\propto P(e, f | D)P(D | c) \end{aligned}$$

Causal and diagnostic support

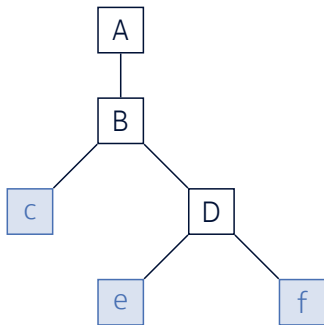


For each X , we compute

- Diagnostic support $P(\text{evidence}_{\text{below } X} \mid X)$
Bottom-up propagation
- Causal support $P(X \mid \text{evidence}_{\text{above } X})$
Top-down propagation

$$\text{BEL}(X) \propto P(\text{evidence}_{\text{below } X} \mid X) \times P(X \mid \text{evidence}_{\text{above } X})$$

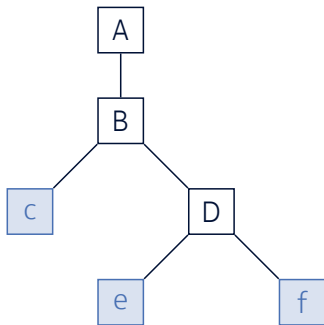
Diagnostic support $P(\text{evidence below } X \mid X)$



Bottom-up propagation



Diagnostic support $P(\text{evidence below } X \mid X)$

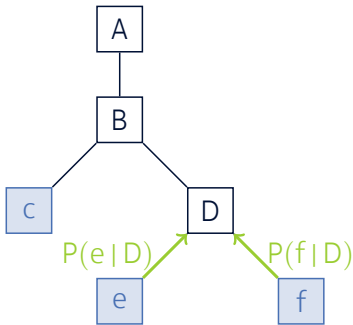


Bottom-up propagation

- $P(e, f \mid D) = P(e \mid D)P(f \mid D)$



Diagnostic support $P(\text{evidence below } X \mid X)$

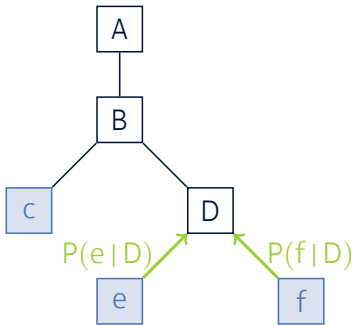


Bottom-up propagation

- $P(e, f | D) = P(e | D)P(f | D)$



Diagnostic support $P(\text{evidence below } X \mid X)$

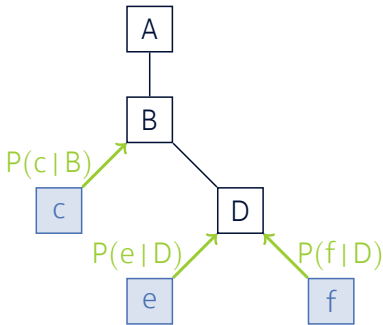


Bottom-up propagation

- $P(e, f \mid D) = P(e \mid D)P(f \mid D)$
- $P(c, e, f \mid B) = P(c \mid B)P(e, f \mid B)$



Diagnostic support $P(\text{evidence below } X \mid X)$

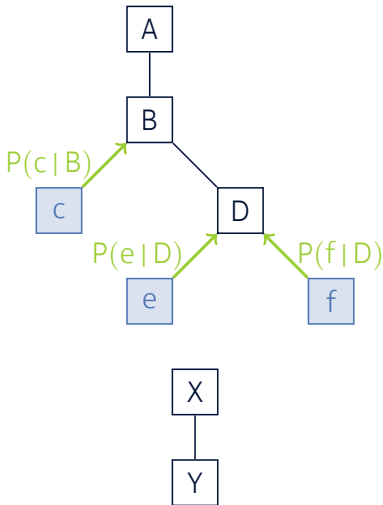


Bottom-up propagation

- $P(e, f | D) = P(e | D)P(f | D)$
- $P(c, e, f | B) = P(c | B)P(e, f | B)$



Diagnostic support $P(\text{evidence below } X \mid X)$



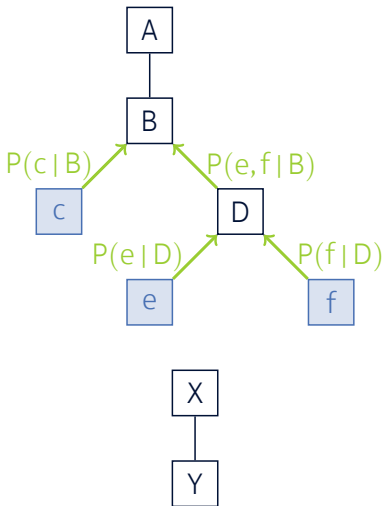
Bottom-up propagation

- $P(e, f | D) = P(e | D)P(f | D)$
- $P(c, e, f | B) = P(c | B)P(e, f | B)$

Compute $P(e, f | B)$:

$$\begin{aligned} P(e, f | B) &= \sum_D P(e, f | B, D)P(D | B) \\ &= \sum_D P(e, f | D)P(D | B) \end{aligned}$$

Diagnostic support $P(\text{evidence below } X \mid X)$



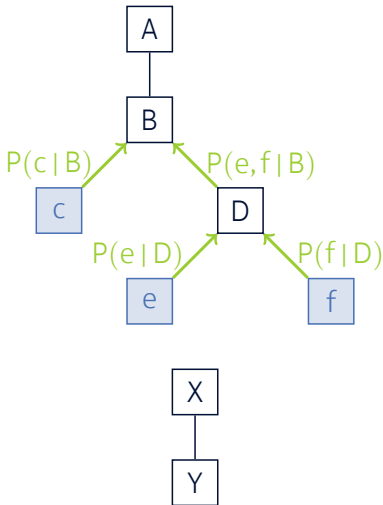
Bottom-up propagation

- $P(e, f | D) = P(e | D)P(f | D)$
- $P(c, e, f | B) = P(c | B)P(e, f | B)$

Compute $P(e, f | B)$:

$$\begin{aligned} P(e, f | B) &= \sum_D P(e, f | B, D)P(D | B) \\ &= \sum_D P(e, f | D)P(D | B) \end{aligned}$$

Diagnostic support $P(\text{evidence below } X \mid X)$

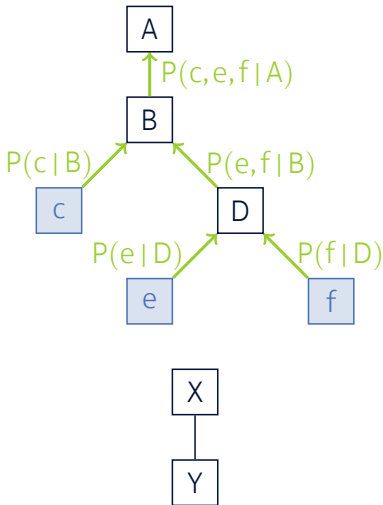


- $P(c, e, f | A) = P(c, e, f | A)$

Compute $P(c, e, f | A)$:

$$\begin{aligned} P(c, e, f | A) &= \sum_B P(c, e, f | A, B) P(B | A) \\ &= \sum_B P(c, e, f | B) P(B | A) \end{aligned}$$

Diagnostic support $P(\text{evidence below } X \mid X)$

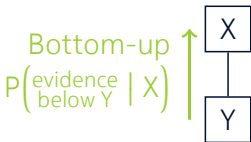
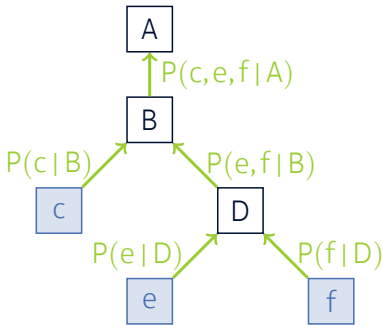


- $P(c, e, f | A) = P(c, e, f | A)$

Compute $P(c, e, f | A)$:

$$\begin{aligned} P(c, e, f | A) &= \sum_B P(c, e, f | A, B) P(B | A) \\ &= \sum_B P(c, e, f | B) P(B | A) \end{aligned}$$

Diagnostic support $P(\text{evidence below } X \mid X)$

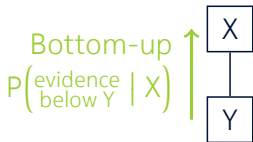
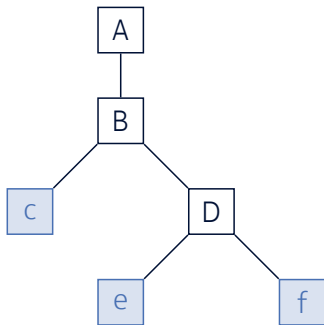


- $P(c, e, f | A) = P(c, e, f | A)$

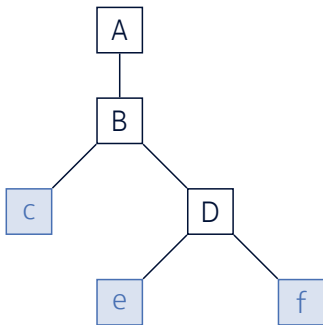
Compute $P(c, e, f | A)$:

$$\begin{aligned} P(c, e, f | A) &= \sum_B P(c, e, f | A, B) P(B | A) \\ &= \sum_B P(c, e, f | B) P(B | A) \end{aligned}$$

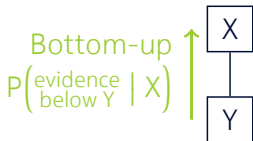
Causal support $P(X \mid \text{evidence above } X)$



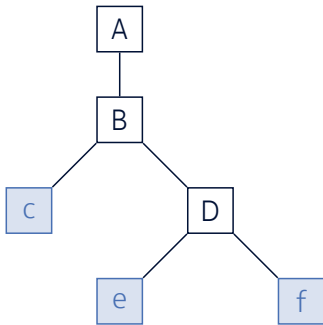
Causal support $P(X \mid \text{evidence above } X)$



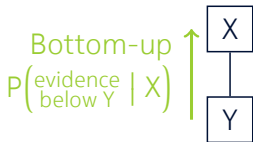
- $P(A)$: parameter



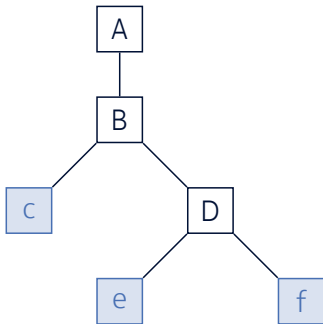
Causal support $P(X \mid \text{evidence above } X)$



- $P(A)$: parameter
→ $BEL(A) \propto P(c, e, f \mid A)P(A)$



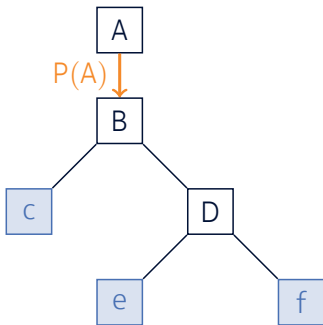
Causal support $P(X \mid \text{evidence above } X)$



- $P(A)$: parameter
→ $BEL(A) \propto P(c, e, f | A)P(A)$
- $P(B) = \sum_A P(B | A)P(A)$



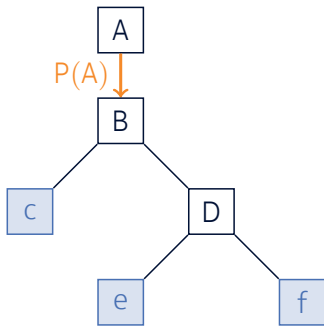
Causal support $P(X \mid \text{evidence above } X)$



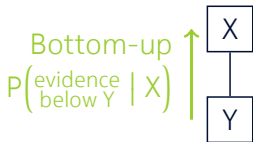
- $P(A)$: parameter
→ $BEL(A) \propto P(c, e, f \mid A)P(A)$
- $P(B) = \sum_A P(B \mid A)P(A)$



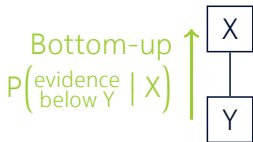
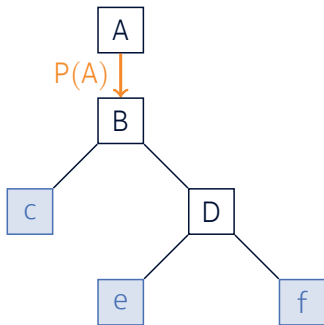
Causal support $P(X \mid \text{evidence above } X)$



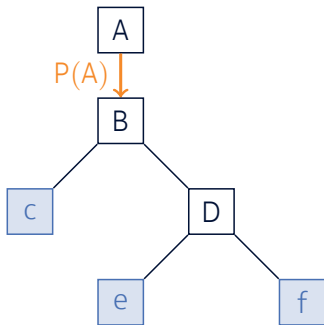
- $P(A)$: parameter
→ $BEL(A) \propto P(c, e, f \mid A)P(A)$
- $P(B) = \sum_A P(B \mid A)P(A)$
→ $BEL(B) \propto P(c, e, f \mid B)P(B)$



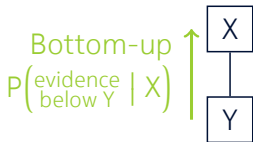
Causal support $P(X \mid \text{evidence above } X)$



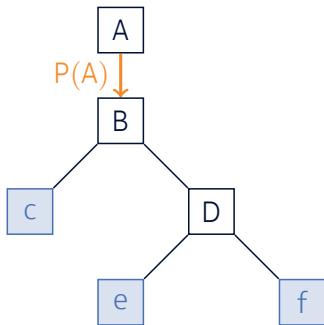
Causal support $P(X \mid \text{evidence above } X)$



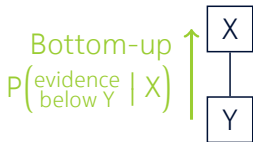
- $P(D \mid c) = \sum_B P(D \mid B, c)P(B \mid c)$



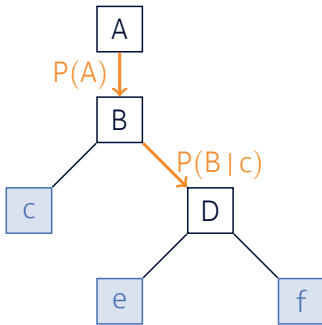
Causal support $P(X \mid \text{evidence above } X)$



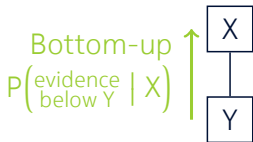
$$\begin{aligned} \bullet P(D \mid c) &= \sum_B P(D \mid B, c) P(B \mid c) \\ &= \sum_B P(D \mid B) P(B \mid c) \end{aligned}$$



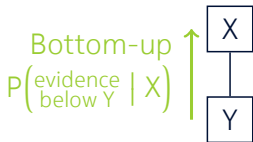
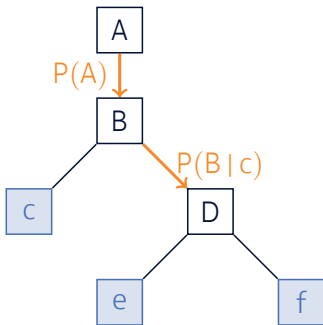
Causal support $P(X \mid \text{evidence above } X)$



- $$P(D|c) = \sum_B P(D|B,c)P(B|c)$$
$$= \sum_B P(D|B)P(B|c)$$



Causal support $P(X \mid \text{evidence above } X)$



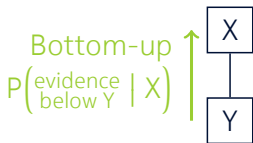
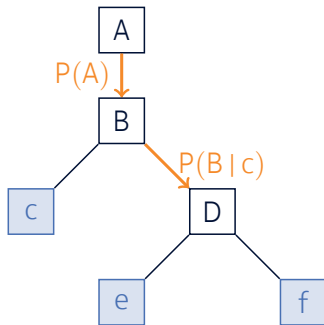
- $$P(D | c) = \sum_B P(D | B, c) P(B | c)$$
$$= \sum_B P(D | B) P(B | c)$$

Compute $P(B | c)$:

$$P(B | c, e, f) = \frac{P(e, f | B, c) P(B | c)}{P(e, f | c)}$$

$$\text{i.e. } P(B | c) \propto \frac{BEL(B)}{P(e, f | B)}$$

Causal support $P(X \mid \text{evidence above } X)$



- $$P(D \mid c) = \sum_B P(D \mid B, c) P(B \mid c)$$
$$= \sum_B P(D \mid B) P(B \mid c)$$

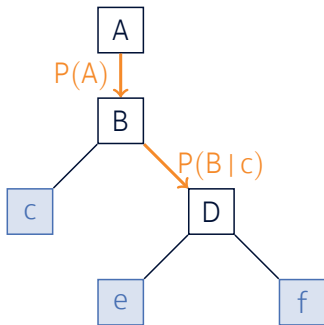
Compute $P(B \mid c)$:

$$P(B \mid c, e, f) = \frac{P(e, f \mid B, c) P(B \mid c)}{P(e, f \mid c)}$$

$$\text{i.e. } P(B \mid c) \propto \frac{\text{BEL}(B)}{P(e, f \mid B)}$$

$$\rightarrow \text{BEL}(D) \propto P(e, f \mid D) P(D \mid c)$$

Causal support $P(X \mid \text{evidence above } X)$



- $$P(D|c) = \sum_B P(D|B,c)P(B|c)$$

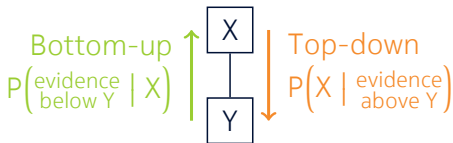
$$= \sum_B P(D|B)P(B|c)$$

Compute $P(B|c)$:

$$P(B|c,e,f) = \frac{P(e,f|B,c)P(B|c)}{P(e,f|c)}$$

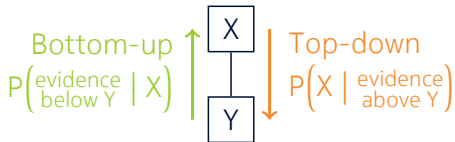
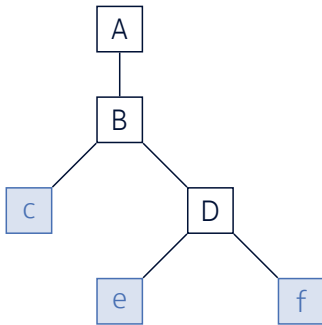
i.e. $P(B|c) \propto \frac{BEL(B)}{P(e,f|B)}$

$\rightarrow BEL(D) \propto P(e,f|D)P(D|c)$

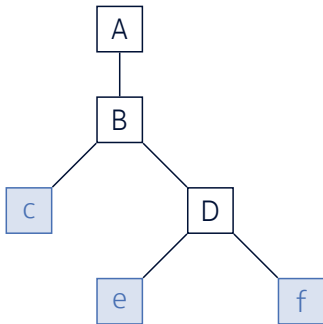


Summary

Algorithm

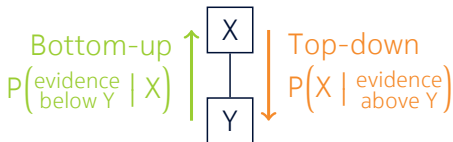


Summary

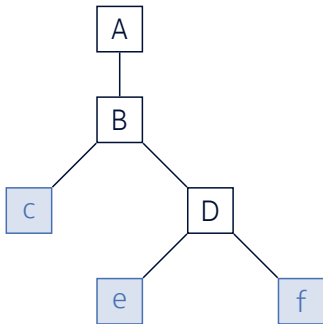


Algorithm

- Diagnostic support $P(\text{evidence}_{\text{below } X} \mid X)$
Bottom-up propagation

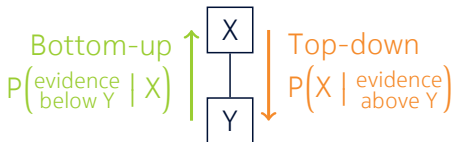


Summary

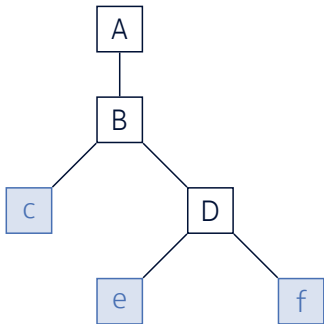


Algorithm

- Diagnostic support $P(\text{evidence below } X \mid X)$
Bottom-up propagation
- Causal support $P(X \mid \text{evidence above } X)$
Top-down propagation



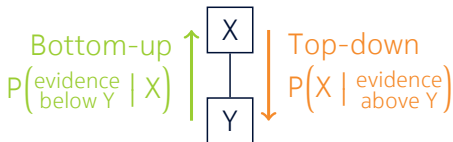
Summary



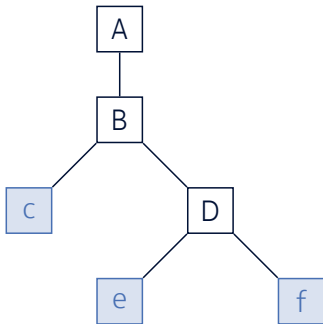
Algorithm

- Diagnostic support $P(\text{evidence below } X \mid X)$
Bottom-up propagation
- Causal support $P(X \mid \text{evidence above } X)$
Top-down propagation

In general



Summary

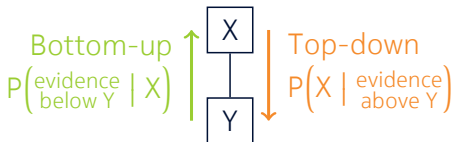


Algorithm

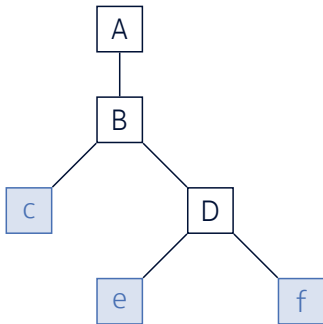
- Diagnostic support $P(\text{evidence below } X \mid X)$
Bottom-up propagation
- Causal support $P(X \mid \text{evidence above } X)$
Top-down propagation

In general

- Use a topological ordering



Summary

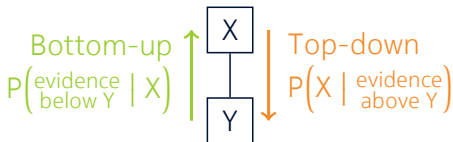


Algorithm

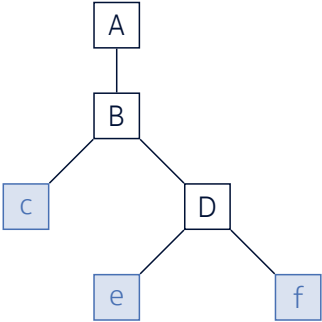
- Diagnostic support $P(\text{evidence below } X \mid X)$
Bottom-up propagation
- Causal support $P(X \mid \text{evidence above } X)$
Top-down propagation

In general

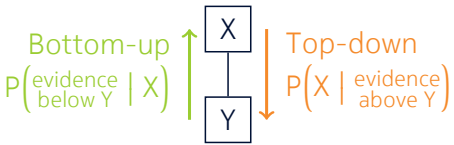
- Use a topological ordering
- Complexity: $O(rd^\downarrow + r^2 + r)$



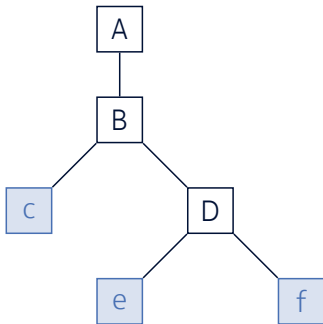
Additional remarks



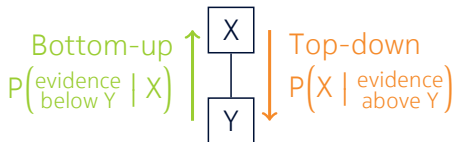
- If the evidence node is not a leaf: add a phantom node



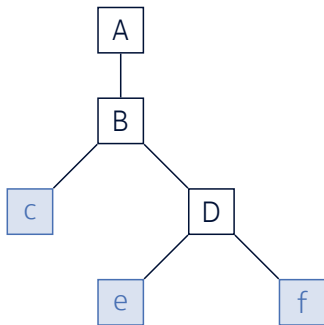
Additional remarks



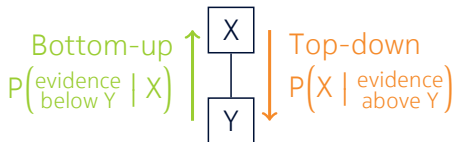
- If the evidence node is not a leaf: add a phantom node
- The calculations can be written as matrix products
 - Belief, causal and diagnostic supports, messages ~ Vectors
 - CPT ~ Matrix



Additional remarks



- If the evidence node is not a leaf: add a phantom node
- The calculations can be written as matrix products
 - Belief, causal and diagnostic supports, messages ~ Vectors
 - CPT ~ Matrix



- Asynchronous / parallel updates: acknowledgements (Pearl, 1982, 1986, 1988)

Outline

Reminders on probability theory

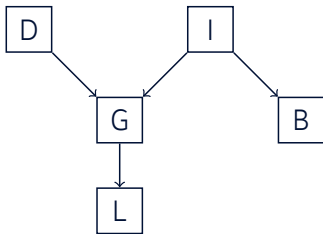
Bayesian networks

Belief propagation in trees

Belief propagation in polytrees

Polytree (or singly-connected) Bayesian network

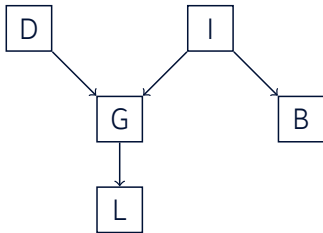
The underlying undirected graph is a tree



Polytree (or singly-connected) Bayesian network

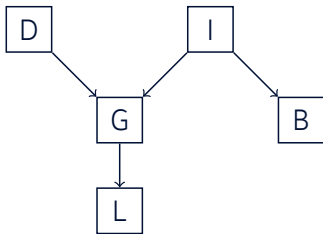
The underlying undirected graph is a tree

Separation properties



Polytree (or singly-connected) Bayesian network

The underlying undirected graph is a tree

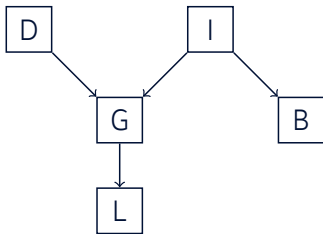


Separation properties

- Given a node, the non-descendants and the subtrees rooted at each child are independent

Polytree (or singly-connected) Bayesian network

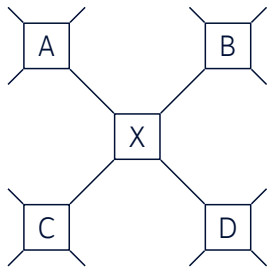
The underlying undirected graph is a tree



Separation properties

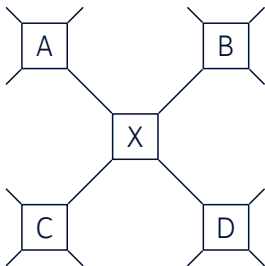
- Given a node, the non-descendants and the subtrees rooted at each child are independent
- If we don't condition on a node nor any of its descendants, the inversed subtrees rooted at its ancestors are independent

Belief propagation



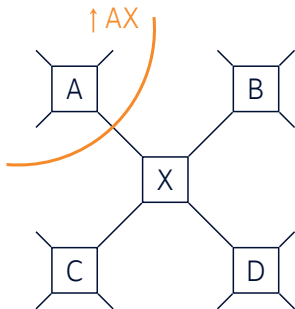
Belief propagation

$$P(X | \text{evidence}) \propto P(\text{evidence}_{\text{below } X} | X) \times P(X | \text{evidence}_{\text{above } X})$$



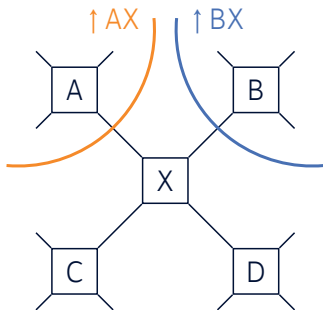
Belief propagation

$$P(X | \text{evidence}) \propto P(\text{evidence below } X | X) \times P(X | \text{evidence above } X)$$



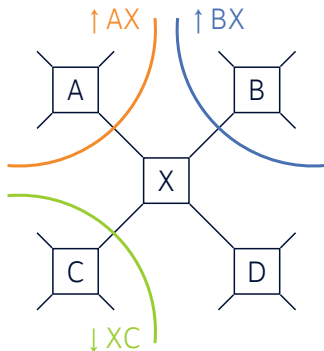
Belief propagation

$$P(X | \text{evidence}) \propto P(\text{evidence below } X | X) \times P(X | \text{evidence above } X)$$



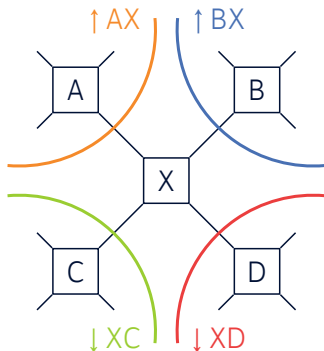
Belief propagation

$$P(X | \text{evidence}) \propto P(\text{evidence below } X | X) \times P(X | \text{evidence above } X)$$



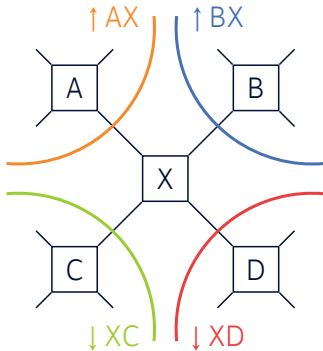
Belief propagation

$$P(X | \text{evidence}) \propto P(\text{evidence below } X | X) \times P(X | \text{evidence above } X)$$



Belief propagation

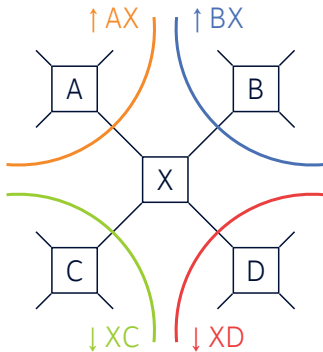
$$P(X | \text{evidence}) \propto P(\text{evidence}_{\text{below } X} | X) \times P(X | \text{evidence}_{\text{above } X})$$



- Diagnostic support $P(\text{evidence}_{\text{below } X} | X)$
Bottom-up propagation
↓ XC and ↓ XD are independent given X

Belief propagation

$$P(X | \text{evidence}) \propto P(\text{evidence}_{\text{below } X} | X) \times P(X | \text{evidence}_{\text{above } X})$$



- Diagnostic support $P(\text{evidence}_{\text{below } X} | X)$
Bottom-up propagation
 $\downarrow XC$ and $\downarrow XD$ are independent given X
- Causal support $P(X | \text{evidence}_{\text{above } X})$
Top-down propagation
 $\uparrow AX$ and $\uparrow BX$ are independent

Conclusion

Conclusion

- **Bayesian networks**

A memory-efficient way of storing a PMF by leveraging conditional independencies between variables

Conclusion

- **Bayesian networks**

A memory-efficient way of storing a PMF by leveraging conditional independencies between variables

- **Belief propagation**

A time-efficient algorithm for computing the belief

- Asynchronous, parallelizable
- Exact in (poly)trees
- In general, extended to the junction tree algorithm and to other (approximate) algorithms