

Fundamentals of AI and KR - Module 3

2. Bayesian network representation

Paolo Torroni

Fall 2023

Notice

Credits

The present slides are largely an adaptation of existing material, including:

- slides from [Russel & Norvig](#)
- slides by [Daphne Koller](#) on Probabilistic Graphical Models
- slides by Fabrizio Riguzzi on [Data Mining and Analytics](#)

I am especially grateful to these authors.

Downloading and sharing

A copy of these slides can be downloaded from [virtuale](#) and stored for personal use only. Please do not redistribute.

Table of Contents

- Independence
- Bayesian network representation

Independence



Independence

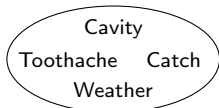
A and B are independent, denoted $\mathbf{P} \models (A \perp B)$, iff

$$\mathbf{P}(A|B) = \mathbf{P}(A) \quad \text{or} \quad \mathbf{P}(B|A) = \mathbf{P}(B) \quad \text{or} \quad \mathbf{P}(A, B) = \mathbf{P}(A)\mathbf{P}(B)$$

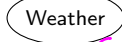
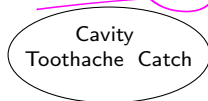
1

2

3



decomposes into



$$\mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather}) \\ = \mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity})\mathbf{P}(\text{Weather})$$

$$\mathbf{P} \models (W \perp T, \text{Cat}, \text{Cav})$$

32 entries reduced to 12; for n independent biased coins, $2^n \rightarrow n$

Absolute (marginal) independence powerful but rare.

Dentistry is a large field with hundreds of variables, none of which are independent. What to do?



Conditional independence

$P(\textit{Toothache}, \textit{Cavity}, \textit{Catch})$ has $2^3 - 1 = 7$ independent entries
 If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:

$$(1) \underline{P(\textit{catch}|\textit{toothache}, \textit{cavity})} = \underline{P(\textit{catch}|\textit{cavity})}$$

The same independence holds if I haven't got a cavity:

$$(2) P(\textit{catch}|\textit{toothache}, \neg \textit{cavity}) = P(\textit{catch}|\neg \textit{cavity})$$

Catch is conditionally independent of *Toothache* given *Cavity*:

$$P(\textit{Catch}|\textit{Toothache}, \textit{Cavity}) = P(\textit{Catch}|\textit{Cavity})$$

Equivalent statements:

- $P(\textit{Toothache}|\textit{Catch}, \textit{Cavity}) = P(\textit{Toothache}|\textit{Cavity})$
- $P(\textit{Toothache}, \textit{Catch}|\textit{Cavity}) = P(\textit{Toothache}|\textit{Cavity})P(\textit{Catch}|\textit{Cavity})$

Notation: $\mathbf{P} \models (\textit{Catch} \perp \textit{Toothache} | \textit{Cavity})$



Conditional independence

Write out full joint distribution using chain rule:

$$\begin{aligned}
 &P(\text{Toothache}, \text{Catch}, \text{Cavity}) \quad \rightarrow \text{7 indep. values} \\
 &= P(\text{Toothache} | \text{Catch}, \text{Cavity}) P(\text{Catch} | \text{Cavity}) P(\text{Cavity}) \\
 &= P(\text{Toothache} | \text{Cavity}) P(\text{Catch} | \text{Cavity}) P(\text{Cavity})
 \end{aligned}$$

(Note: In the original image, a bracket under the first two terms of the second equation is labeled with a '2', indicating 2 independent values for that part.)

i.e., $2 + 2 + 1 = 5$ independent numbers

In most cases, the use of conditional independence reduces the size of the representation of the joint distribution from exponential in n to linear in n .

Conditional independence is our most basic and robust form of knowledge about uncertain environments.



Bayes' Rule

Product rule $P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$

$$\Rightarrow \text{Bayes' rule } P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

or in distribution form

$$\rightarrow P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \alpha P(X|Y)P(Y)$$

Useful for assessing diagnostic probability from causal probability:

$$\rightarrow P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$



Example of diagnosis using Bayes' Rule

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$

Handwritten annotations: Purple arrows point from the text below to the terms in the equation. One arrow points from 'Cause' to 'P(Cause)', another from 'Effect' to 'P(Effect|Cause)', and a third from 'Effect' to 'P(Effect)'.

Say 1 individual in 50,000 suffers from meningitis, 1% from a stiff neck, and 70% of the times meningitis causes a stiff neck. *What is the probability that an individual with a stiff neck has meningitis?*



Example of diagnosis using Bayes' Rule

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$

Say 1 individual in 50,000 suffers from meningitis, 1% from a stiff neck, and 70% of the times meningitis causes a stiff neck. *What is the probability that an individual with a stiff neck has meningitis?*

Let M be meningitis and S be stiff neck.

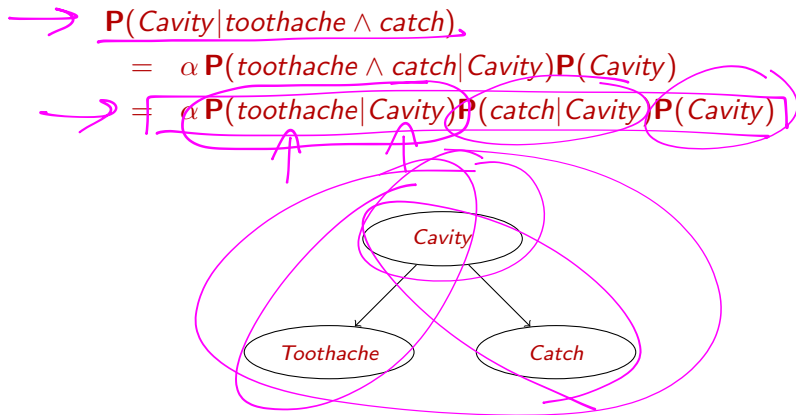
$P(m) = 1/50,000$, $P(s) = 0.01$, $P(s|m) = 0.7$.

$$P(m|s) = \frac{P(s|m)P(m)}{P(s)} = \frac{0.7 \times 1/50,000}{0.01} = 0.0014$$

Note: posterior probability of meningitis still very small!

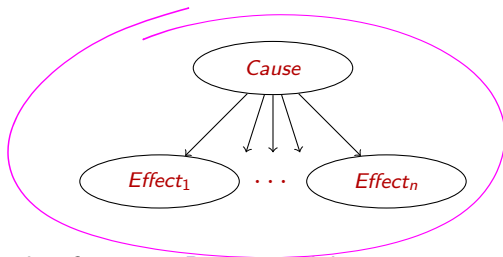


Bayes' Rule and conditional independence





Bayes' Rule and conditional independence



This is an example of a **naive Bayes** model:

$$\underline{\mathbf{P}(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n)} = \mathbf{P}(\text{Cause}) \prod_i \mathbf{P}(\text{Effect}_i | \text{Cause})$$

Total number of parameters is linear in n



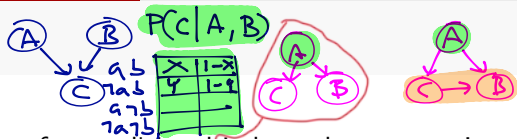
Summary so far

- Probability is a rigorous formalism for uncertain knowledge
- Joint probability distribution specifies probability of every atomic event
- Queries can be answered by summing over atomic events
- For nontrivial domains, we must find a way to reduce the joint size
- Independence and conditional independence provide the tools

Bayesian network representation



Bayesian networks



A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions.

Syntax:

known A, C and B are independent



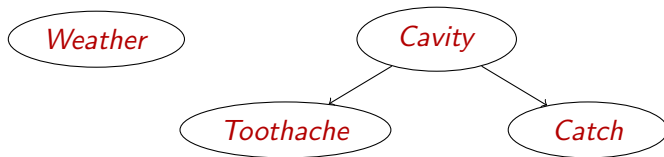
- a set of nodes, one per variable
- a directed, acyclic graph (link \approx “directly influences”)
- a conditional distribution for each node given its parents:
 $P(X_i | Parents(X_i))$

In the simplest case, conditional distribution represented as a conditional probability table (CPT) giving the distribution over X_i for each combination of parent values



Example

Topology of network encodes conditional independence assertions:



- *Weather* is independent of the other variables
- *Toothache* and *Catch* are conditionally independent given *Cavity*



Example

I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?

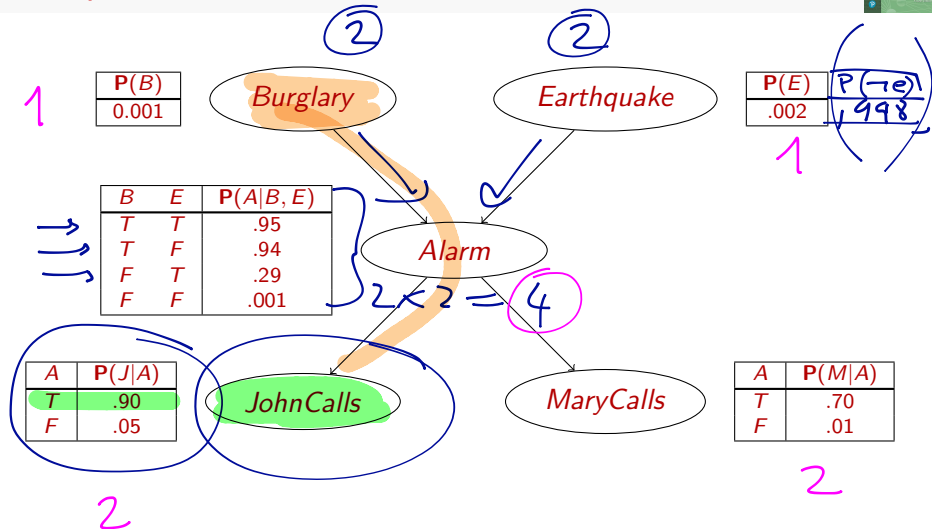
Variables: Burglar, Earthquake, Alarm, JohnCalls, MaryCalls

Network topology reflects "causal" knowledge:

- A burglar can set the alarm off
- An earthquake can set the alarm off
- The alarm can cause Mary to call
- The alarm can cause John to call



Example

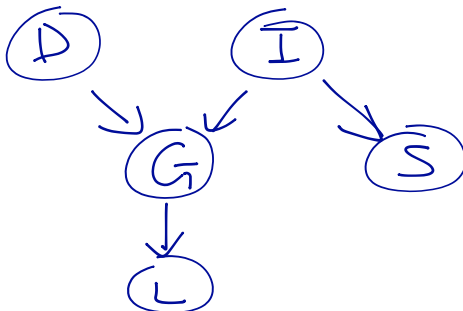


Reasoning patterns



The student network

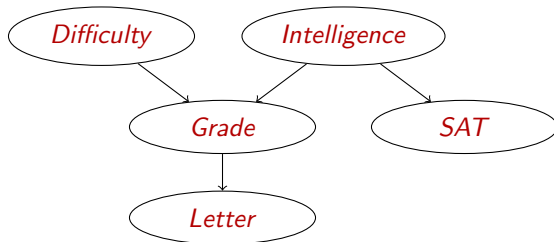
A student's **grade** depends on **intelligence** and on the **difficulty** of the course. **SAT** scores are correlated with intelligence. A professor writes recommendation **letters** by only looking at grades.



Reasoning patterns

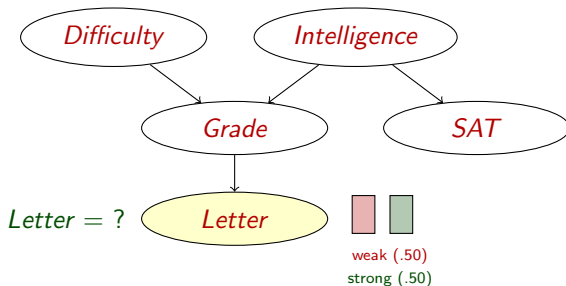
The student network

A student's **grade** depends on **intelligence** and on the **difficulty** of the course. **SAT** scores are correlated with intelligence. A professor writes recommendation **letters** by only looking at grades.



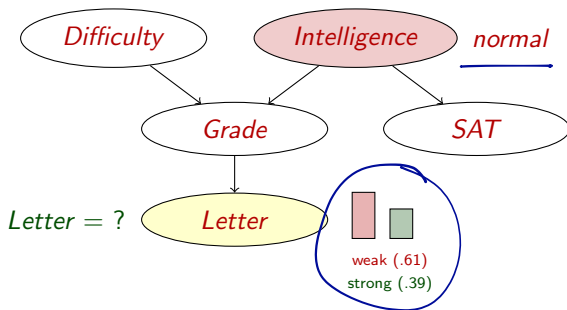
Reasoning patterns

- Causal: will George get a strong reference letter? (prediction)



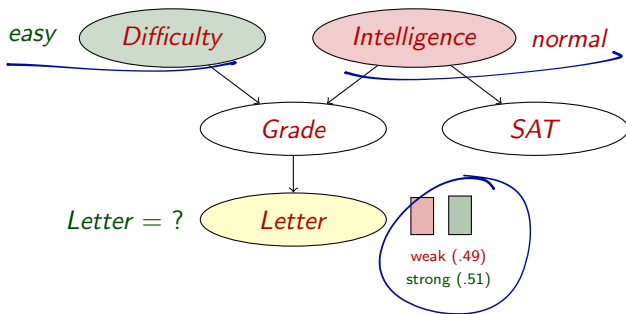
Reasoning patterns

- Causal: will George get a strong reference letter? (prediction)



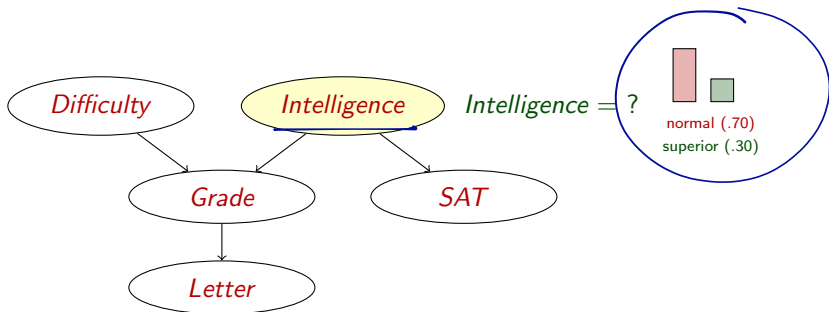
Reasoning patterns

- Causal: will George get a strong reference letter? (prediction)



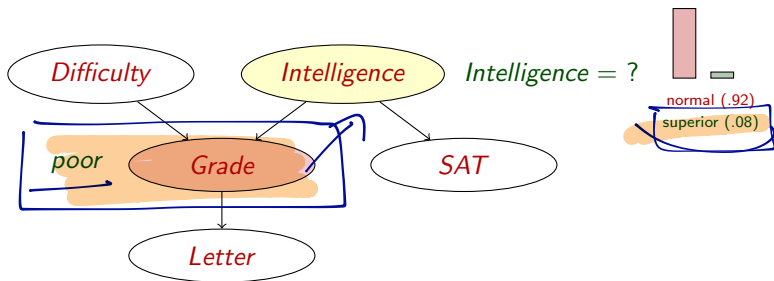
Reasoning patterns

- **Causal**: will George get a strong reference letter? (**prediction**)
- **Evidential**: is George a good potential recruit? (**explanation**)



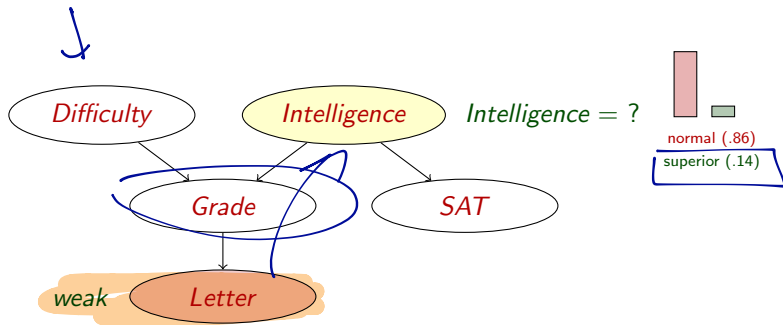
Reasoning patterns

- **Causal**: will George get a strong reference letter? (prediction)
- **Evidential**: is George a good potential recruit? (explanation)



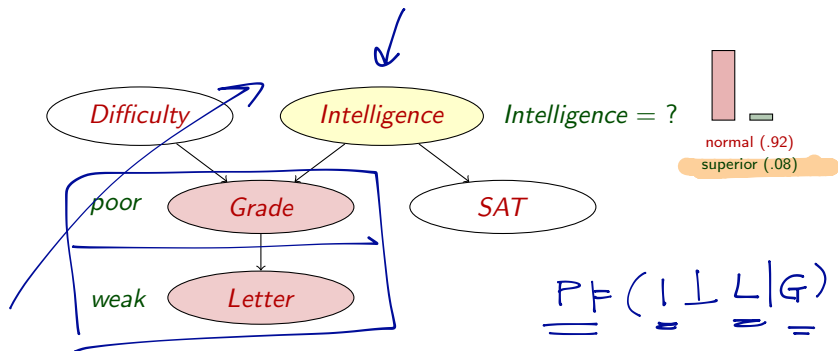
Reasoning patterns

- **Causal**: will George get a strong reference letter? (**prediction**)
- **Evidential**: is George a good potential recruit? (**explanation**)



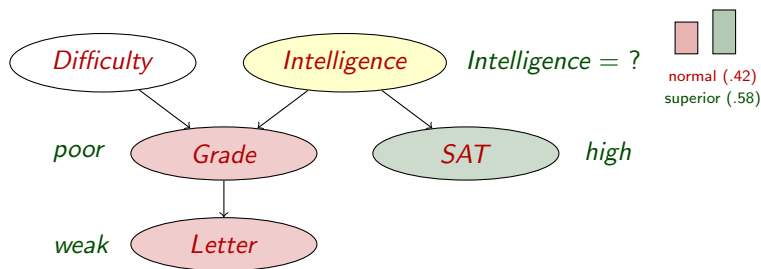
Reasoning patterns

- **Causal**: will George get a strong reference letter? (**prediction**)
- **Evidential**: is George a good potential recruit? (**explanation**)



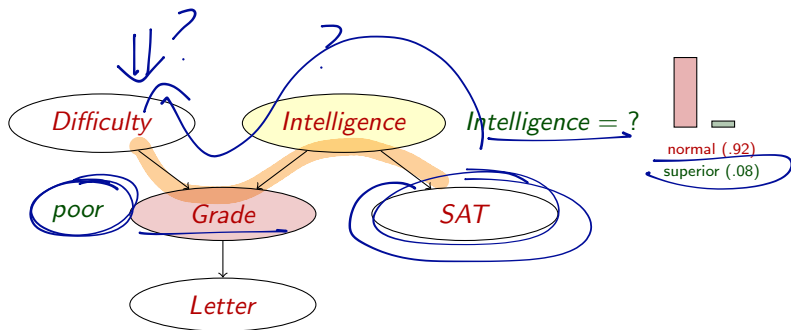
Reasoning patterns

- **Causal**: will George get a strong reference letter? (**prediction**)
- **Evidential**: is George a good potential recruit? (**explanation**)



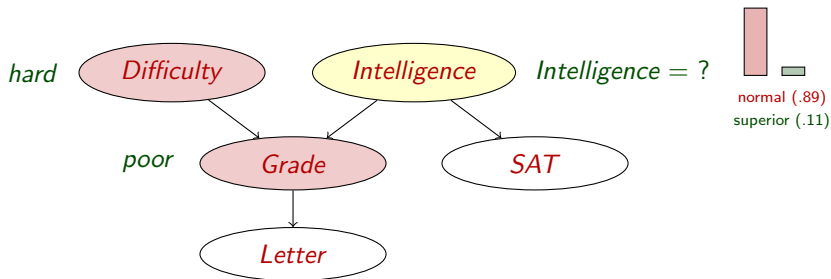
Reasoning patterns

- **Causal**: will George get a strong reference letter? (prediction)
- **Evidential**: is George a good potential recruit? (explanation)
- **Intercausal**: why did George score low/high? (explaining away)



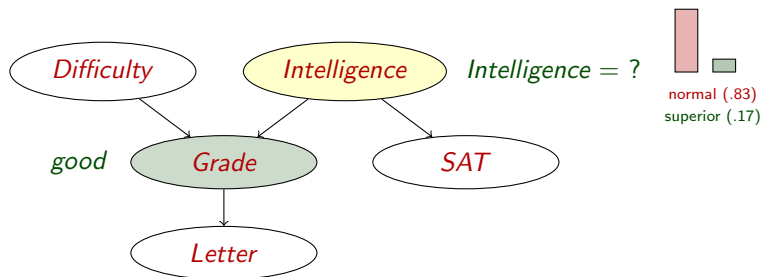
Reasoning patterns

- **Causal**: will George get a strong reference letter? (**prediction**)
- **Evidential**: is George a good potential recruit? (**explanation**)
- **Intercausal**: why did George score low/high? (**explaining away**)



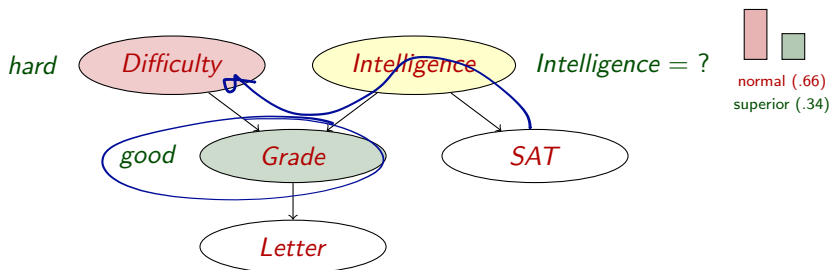
Reasoning patterns

- **Causal**: will George get a strong reference letter? (**prediction**)
- **Evidential**: is George a good potential recruit? (**explanation**)
- **Intercausal**: why did George score low/high? (**explaining away**)



Reasoning patterns

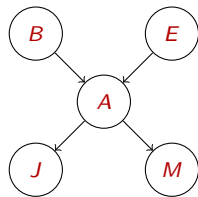
- **Causal**: will George get a strong reference letter? (**prediction**)
- **Evidential**: is George a good potential recruit? (**explanation**)
- **Intercausal**: why did George score low/high? (**explaining away**)





Compactness

A CPT for Boolean X_i with k Boolean parents has 2^k rows for the combinations of parent values
 Each row requires one number p for $X_i = \text{true}$
 (the number for $X_i = \text{false}$ is just $1 - p$)



If each variable has no more than k parents, the complete network requires $O(n \cdot 2^k)$ numbers

I.e., grows linearly with n , vs. $O(2^n)$ for the full joint distribution

For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers (vs. $2^5 - 1 = 31$)

For student net, $1 + 1 + 8 + 2 + 3 = 15$ numbers (vs. $2^4 \times 3 - 1 = 47$)



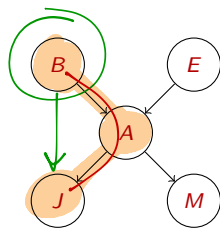
Global semantics

Global semantics defines the full joint distribution as the product of the local conditional distributions:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e) \rightarrow$

$$\rightarrow \frac{P(j|a)P(m|a)P(a|\neg b, e)P(\neg b)P(e)}{P(j|a, m, \neg b, e)P(a|m, \neg b, e)P(\dots)}$$





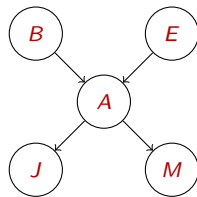
Global semantics

Global semantics defines the full joint distribution as the product of the local conditional distributions:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

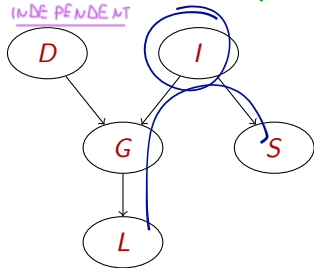
e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$\begin{aligned} &= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) \\ &= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\ &\approx 0.00063 \end{aligned}$$



Basic independencies in the Student network

IF YOU DON'T KNOW
THE GRADE, POINTS
AND DIFFICULTY
ARE INDEPENDENT



$$P \models (S \perp D)$$

$$P \models (S \perp D, L, G | I)$$

What independencies?

- $P \models (L \perp \dots ?)$

- $P \models (S \perp \dots ?)$

- $P \models (G \perp \dots ?)$

- $P \models (I \perp \dots ?)$

- $P \models (D \perp \dots ?)$

$$P \models (S \perp L | G)$$

$$P \models (S \perp D | G)$$

$$P \models (L \perp D, I, S | G)$$

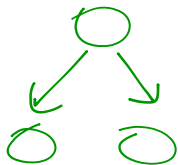
$$P \models (\underbrace{L \perp S}_{(S \perp L | I)} | I)$$

IF YOU KNOW THE GRADE,
POINTS AND DIFFICULTY
ARE DEPENDENT

Basic independencies in the Student network

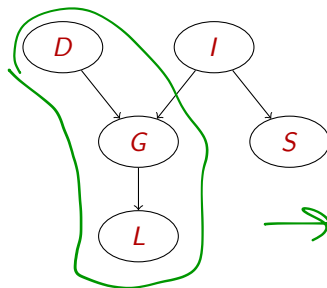
What independencies?

- $\mathbf{P} \models (L \perp I, D, S | G)$
- $\mathbf{P} \models (S \perp G, D, L | I)$
- $\mathbf{P} \models (G \perp S | I)$
- $\mathbf{P} \models (I \perp D)$
- $\mathbf{P} \models (D \perp I, S)$
- ...



V-structure

Flow of probabilistic influence



Or else, when could X influence Y ?

- $X \rightarrow Y$ (direct cause)
- $X \leftarrow Y$ (direct effect)
- $X \rightarrow Z \rightarrow Y$ (causal trail)
- $X \leftarrow Z \leftarrow Y$ (evidential trail)
- $X \leftarrow Z \rightarrow Y$ (common cause)
- $X \rightarrow Z \leftarrow Y$ (common effect)

Definition (active two-edge trail)

If influence can flow from X to Y via Z , the trail $X \rightleftharpoons Z \rightleftharpoons Y$ is **active**

Flow of probabilistic influence: active trails

Consider a longer trail $X_1 \Rightarrow \dots \Rightarrow X_n$.

For influence to flow from X_1 to X_n , it needs to flow through every single node on the trail

This is true if and only if every two-edge trail $X_{i-1} \Rightarrow X_i \Rightarrow X_{i+1}$ along the trail allows influence to flow

Definition (active trail)

Let \mathbf{Z} be a subset of observed variables.

The trail $X_{i-1} \Rightarrow X_i \Rightarrow X_{i+1}$ is active given \mathbf{Z} if

- $\forall X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$, X_i or one of its descendants are in \mathbf{Z}
- no other node along the trail is in \mathbf{Z}

Flow of probabilistic influence: direct separation

Definition (d-separation)

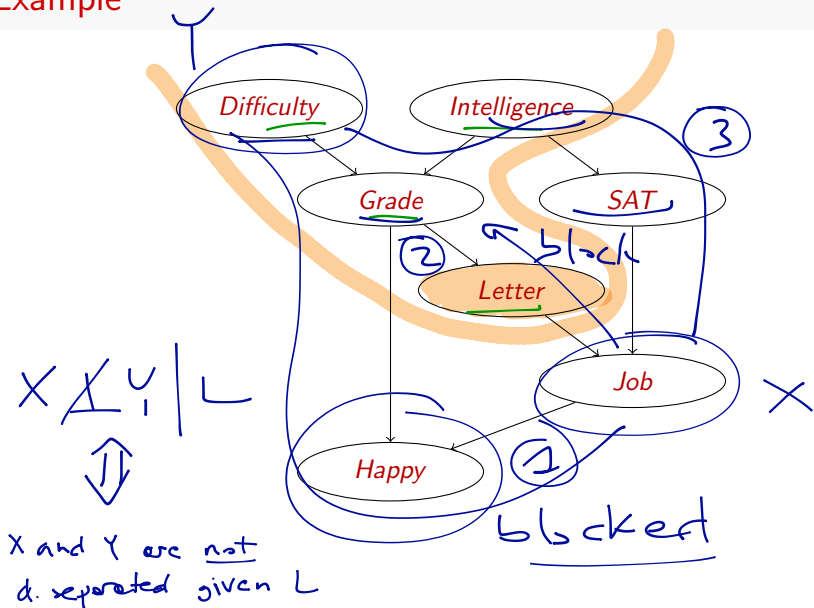
Two sets of nodes \mathbf{X} , \mathbf{Y} are **d-separated** given \mathbf{Z} if there is no active trail between any $X \in \mathbf{X}$ and $Y \in \mathbf{Y}$ given \mathbf{Z}

To determine if \mathbf{X} and \mathbf{Y} are **independent** given \mathbf{Z} :

- ① traverse the graph bottom-up marking all nodes in \mathbf{Z} or having descendants in given \mathbf{Z}
- ② traverse the graph from \mathbf{X} to \mathbf{Y} , stopping if we get to a **blocked** node
- ③ if we can't reach \mathbf{Y} , then \mathbf{X} and \mathbf{Y} are independent

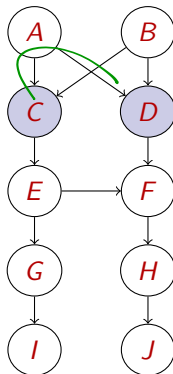
A node is **blocked** if either the middle of an unmarked v-structure, or in \mathbf{Z} (not both)

Example





Example



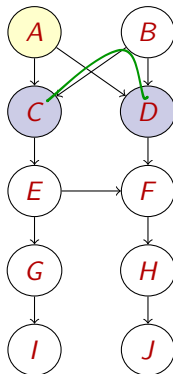
What independences?

- $\mathbf{P} \models (C \perp D)$? no

THERE IS AN ACTIVE TRAIL



Example



What independences?

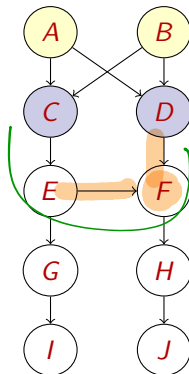
- $\mathbf{P} \models (C \perp D)$?
- $\mathbf{P} \models (C \perp D | A)$? ~



ONE IS ACTIVE BUT THE OTHER
ONE NOT



Example

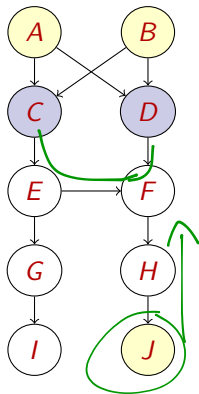


What independences?

- $\mathbf{P} \models (C \perp D)?$
- $\mathbf{P} \models (C \perp D | A)?$
- $\mathbf{P} \models (C \perp D | A, B)?$ *yes*



Example



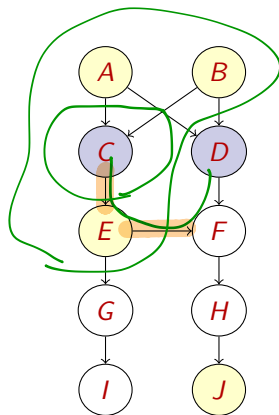
What independences?

- $\mathbf{P} \models (C \perp D)$?
- $\mathbf{P} \models (C \perp D | A)$?
- $\mathbf{P} \models (C \perp D | A, B)$?
- $\mathbf{P} \models (C \perp D | \underline{A, B, J})$?

~



Example



What independences?

- $\mathbf{P} \models (C \perp D)?$
- $\mathbf{P} \models (C \perp D | A)?$
- $\mathbf{P} \models (C \perp D | A, B)?$
- $\mathbf{P} \models (C \perp D | A, B, J)?$
- $\mathbf{P} \models (C \perp D | A, B, E, J)?$ yes

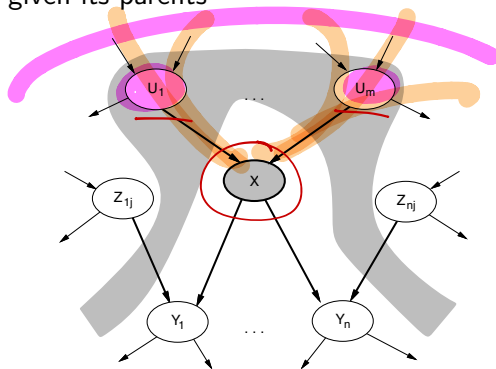
$C \perp$ all other nodes	A, B, E
---------------------------	-----------

↓
"markov blanket"



Local semantics

Local semantics: each node is conditionally independent of its nondescendants given its parents

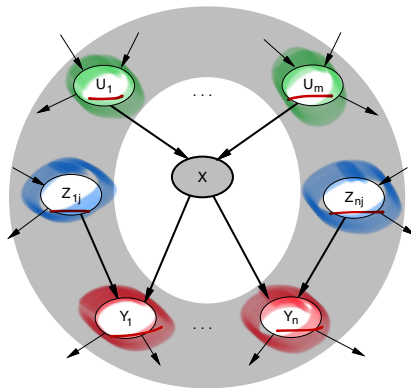


Theorem: Local semantics \Leftrightarrow global semantics



Markov blanket

Each node is conditionally independent of all others given its **Markov blanket**: **parents** + **children** + **children's parents**



Questions?