



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

PROBABILISTIC GRAPHICAL MODEL SESSION # 4 : BAYESIAN MODEL

SEETHA PARAMESWARAN

seetha.p@pilani.bits-pilani.ac.in

The instructor is gratefully acknowledging
the authors who made their course
materials freely available online.

TABLE OF CONTENTS

1 BAYESIAN NETWORK

2 REASONING PATTERNS

3 INDEPENDENCY MAP

BAYESIAN NETWORK

DEFINITION (GLOBAL SEMANTICS)

A Bayesian Network is a directed acyclic graph G whose nodes represent the random variables $\{X_1, X_2, \dots, X_n\}$ and represents a joint distribution via the chain rule for the Bayesian Networks.

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa(X_i)) \quad (1)$$

- Each node is associated with a CPD.

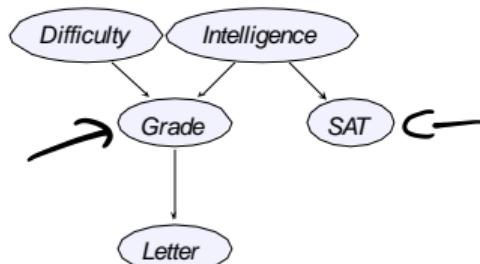


$$CPD(X_i) = P(X_i | Pa(X_i)) \quad \text{and} \quad P(X_1 | \emptyset) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots$$

BAYESIAN NETWORK IS LEGAL

A BN is a legal distribution; if

- $P \geq 0$
 -) P is a product of CPDs.
 -) CPDs are non-negative.
- $\sum P = 1$



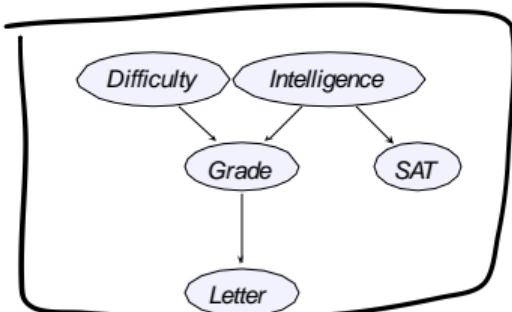
D - I - G - S

~~BAYESIAN NETWORK IS LEGAL~~



A BN is a legal distribution; if

- $P \geq 0$
 -) P is a product of CPDs.
 -) CPDs are non negative.
- $\sum P = 1$



$$\begin{aligned}
 \sum P &= P(I, D, G, S, L) \\
 &= \sum_{D, I, G, S, L} P(I)P(D)P(G|I, D)P(S|I)P(L|G) \\
 &= \underbrace{\sum_{D, I, G, S} P(I)P(D)P(G|I, D)P(S|I)}_{\sum P(L|G)} \\
 &= \underbrace{\sum_{D, I, G} P(I)P(D)P(G|I, D)}_{\sum S P(S|I)} = 1 \\
 &= \underbrace{\sum_{D, I} P(I)P(D)}_{(\alpha + \beta + \gamma)} \underbrace{\sum_G P(G|I, D)}_{(d + e)} \underbrace{\sum_L P(L|G)}_{(\gamma + \delta)} = 1 \\
 &= \left(\sum_I P(I) \right) \left(\sum_D P(D) \right) = 1
 \end{aligned}$$

TABLE OF CONTENTS

1 BAYESIAN NETWORK

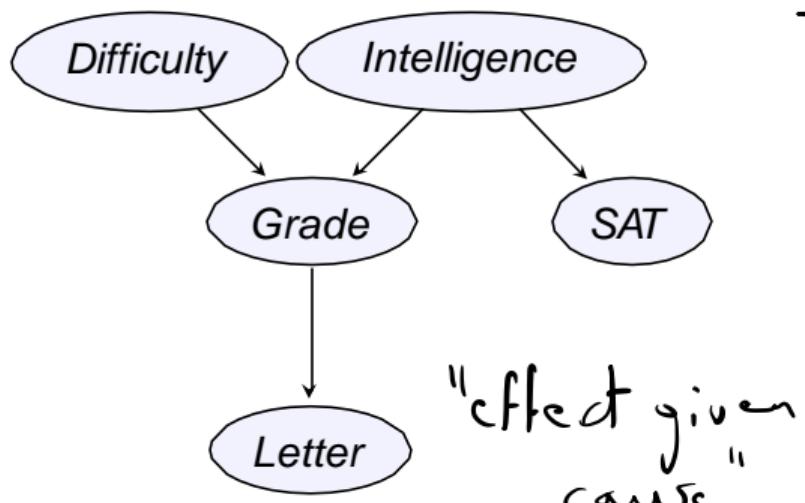
2 REASONING PATTERNS

3 INDEPENDENCY MAP

REASONING PATTERNS

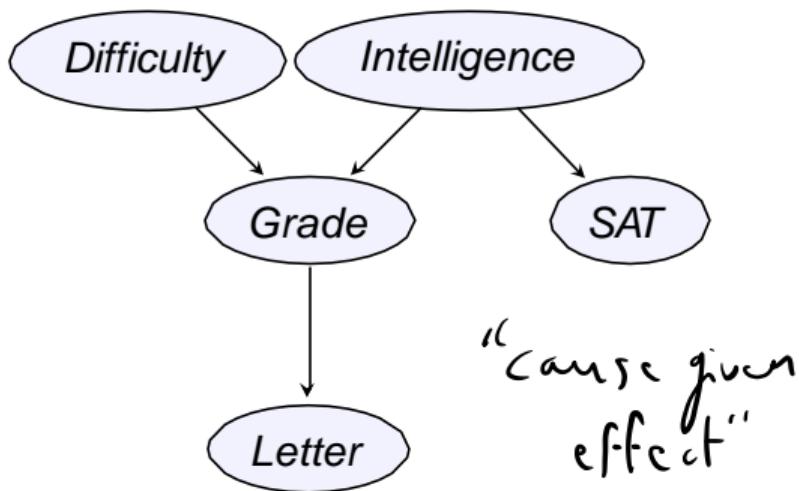
- 1 Causal reasoning
- 2 Evidential reasoning
- 3 Intercausal reasoning

CAUSAL REASONING



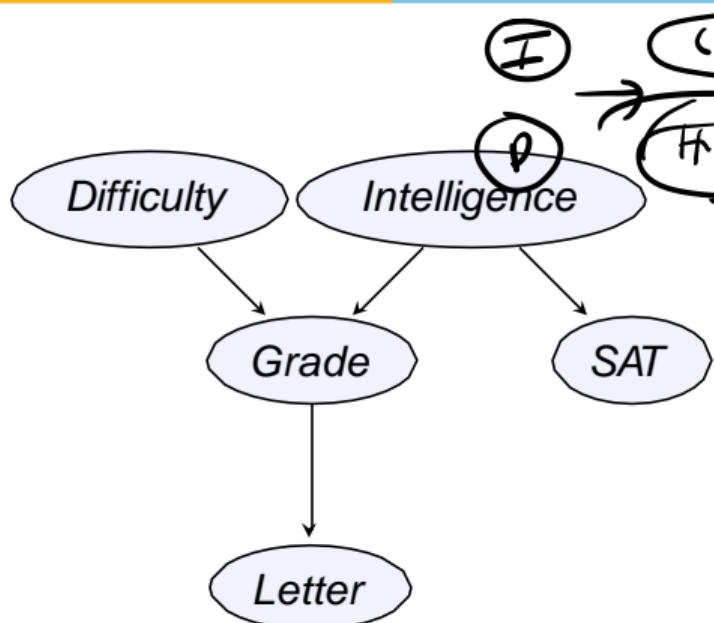
- How likely will a student get a strong recommendation?
 $P(I^1) = ?$
- Given that the student is not so intelligent, what is chance that he gets a strong letter?
 $P(I^1 | I^0) = ?$
- What if the course is easy?
 $P(I^1 | L^0, d^0) = ?$
- Queries that predict the effects of various factors or features are called causal reasoning.

EVIDENTIAL REASONING

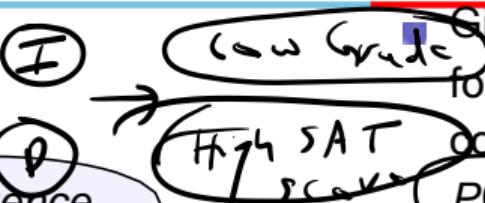


- Given that a student gets C grade for a course, comment on his intelligence.
 $P(i^1|g^3) = ?$
- Given that the student got a weak letter, comment on his intelligence.
 $P(i^1|l^0) = ?$
- $P(i^1|l^0, g^3) = ?$
- Queries that reason from effects to causes are called evidential reasoning.

INTERCAUSAL REASONING



given evidence, prob of multiple causes?



Given that a student gets C grade for a course, and a high SAT score, comment on his intelligence.

$$P(i^1 | g^3, s^1) = ? \text{ could be high}$$

- Does this give any idea regarding the difficulty of the course?

$$P(d^1 | g^3, s^1) = ?$$

- Explained away the poor grade via difficulty of the class. (even though student was intelligent)
- Explaining away is an instance of intercausal reasoning, where different causes of the same effect can interact.

Causal reasoning Example

In the absence of any other information

$$P_{\text{student}}(l') = 0.502$$

How do we calculate this? → From CPDs in
the Bayesian network.

Causal Reasoning Example

From the CPDs we have

$$P(e) = \underbrace{P(e/g)}_{P(g)} P(g) + \underbrace{P(e/g^2)}_{P(g^2)} P(g^2) + \underbrace{P(e/g^3)}_{P(g^3)} P(g^3)$$

$$\begin{aligned} P(g) &= \underbrace{P(g/i^-, d^-)}_{P(i^-, d^-)} P(i^-, d^-) + \underbrace{P(g/i^-, d^+)}_{P(i^-, d^+)} P(i^-, d^+) \\ &\quad + \underbrace{P(g/i^+, d^-)}_{P(i^+, d^-)} P(i^+, d^-) + \underbrace{P(g/i^+, d^+)}_{P(i^+, d^+)} P(i^+, d^+) \end{aligned}$$

Causal Reasoning Example

$$P(i^{\circ}, d^{\circ}) = P(i^{\circ}) P(d^{\circ})$$

We have $P(g) = (0.3)(0.42) + (0.05)(0.28) + (0.9)(0.18)$
 $+ (0.5)(0.12) = \underline{0.362}$

$$P(g^2) = (0.4)(0.42) + (0.25)(0.28) + (0.08)(0.18)
+ (0.3)(0.12) = \underline{0.2884}$$

$$P(g^3) = (0.3)(0.42) + (0.7)(0.28) + (0.02)(0.18) + (0.7)(0.12)
= \underline{0.3496}$$

Causal reasoning Example

$$P(I) = \underline{0.9}(\underline{0.3}C_2) + \underline{0.4}(\underline{0.2884}) + \underline{0.01}(\underline{0.349}) \\ \approx \underline{0.582}$$

If we know that George is not so intelligent, what is the probability that he gets a strong letter of recommendation?

$$P(I'/I_0) \rightarrow ?$$

Causal Reasoning Example

$$P(l^i; i^o) = P(l^i; i^o, d^o) + P(l^i; i^o, \bar{d}) \quad [\text{Marginalization}]$$

$$P(l^i; i^o, d^o) = \frac{P(l^i | g^1) P(g^1 | i^o, d^o) P(i^o, d^o)}{P(l^i | g^1) P(g^1 | i^o, d^o) P(i^o, d^o) + P(l^i | g^2) P(g^2 | i^o, d^o) P(i^o, d^o) + P(l^i | g^3) P(g^3 | i^o, d^o) P(i^o, d^o)}$$

Substituting for the various probabilities we have

$$P(l^i; i^o, d^o) = \underline{0.21546}$$

Causal Reasoning Example

$$\text{Similarly } P(l', i^0, d') = 0.056$$

$$P(l', i^0) = P(l', i^0, d) + P(l', i^0, d') = 0.27146$$

$$P(l' | i^0) = \frac{P(l', i^0)}{P(i^0)} = \frac{0.27146}{0.7} \approx 0.389$$

The probability that the student gets a good letter
of recommendation goes down given low intelligence

Evidential Reasoning Example

Suppose a student received a grade $C \rightarrow g^3$
 what is the probability of high intelligence now?

$$P(i^1 | g^3) = 0.079$$

cause
effect

$$P(d^1) = 0.40 \text{ but } P(d^1 | g^3) = 0.62$$

$$P(i^1 | g^3, s^1) = 0.578$$

"cause given effect"

grade is poor, but SAT score is high, so high intelligence is favoured

Intercausal Reasoning

$$P(i' | g^3) = 0.079$$

If we discover that the subject is hard $\rightarrow d'$
 $P(i' | g^3, d') = 0.11 \rightarrow$ this is a partial explanation
 for the student's low grade (the student can be
intelligent and have still got a poor grade because
 the subject was hard)

Intercausal Reasoning

If the student gets a B grade (g^2) :

$$P(i' | g^2) = 0.175$$

$$P(i' | g^2, d') = 0.34$$

We have explained away the poor grade
using the difficulty of the class.

TABLE OF CONTENTS

1 BAYESIAN NETWORK

2 REASONING PATTERNS

3 INDEPENDENCY MAP

DEPENDENCY IN BN

- A node depends directly only on its parents.
- If the student's grade is known, the quality of his recommendation letter is not influenced by information about any other variable. L is conditionally independent of all other nodes in the network given its parent G .

$$(L \perp \{I, D, S\} \mid G)$$

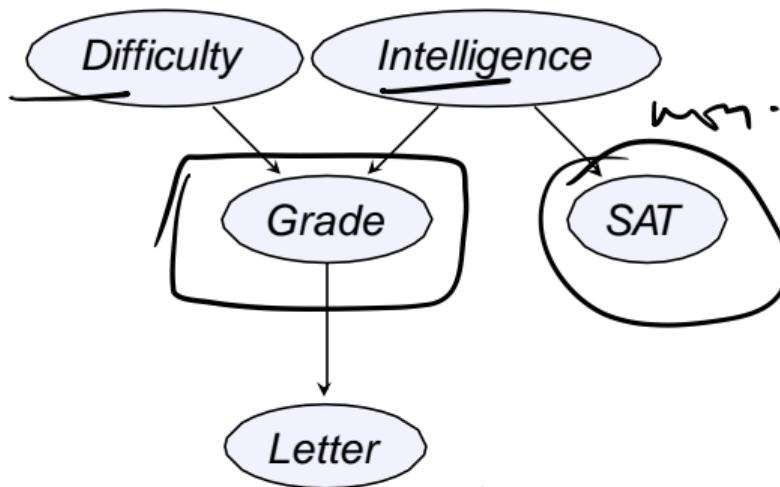
~~$(L \perp I, D \mid G)$~~



- The student's SAT score depends only on his intelligence. S is conditionally independent of all other nodes in the network given its parent I .

$$(S \perp D, G, L \mid I)$$

DEPENDENCY IN BN



- Given the parents, a node can depend on its descendants.

- Does G depend on S , given I and D ?

$$(G \perp S | I, D)$$

- For D , both I and S are non descendants.

$$P(S | I, G, D) = P(S | I) \quad (D \perp I, S)$$

Dependency in BN

Is G independent of C given its parents I and D ?

No. We have:

$$P(g'|i', d', e) > P(g|i', d')$$

If we know that the student got a strong letter of recommendation, it increases the probability that the student had a good grade.

BAYESIAN NETWORK STRUCTURE



DEFINITION (LOCAL SEMANTICS)

A directed acyclic graph G whose nodes represent random variables $\{X_1, X_2, \dots, X_n\}$ and G encodes a set of conditional independence assumptions.

$$\text{For each variable } X_i : (X_i \perp \text{NonDescendants}_{X_i} | \text{Pa}_{X_i}) \quad (2)$$

- Pa_{X_i} represent parents of X_i in G .
- $\text{NonDescendants}_{X_i}$ represent the random variables that are not descendants of X_i .
- $I_L(G)$ represents the set of conditional independence assumptions called local independencies.

BAYESIAN NETWORK SEMANTICS

LOCAL SEMANTICS BN encodes a set of conditional independence assumptions.

For each variable X_i : $(X_i \perp \text{NonDescendants}_{X_i} | \text{Pa}_{X_i})$

GLOBAL SEMANTICS BN represents a joint distribution via the chain rule.

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^{Y_n} P(X_i | \text{Pa}(X_i))$$

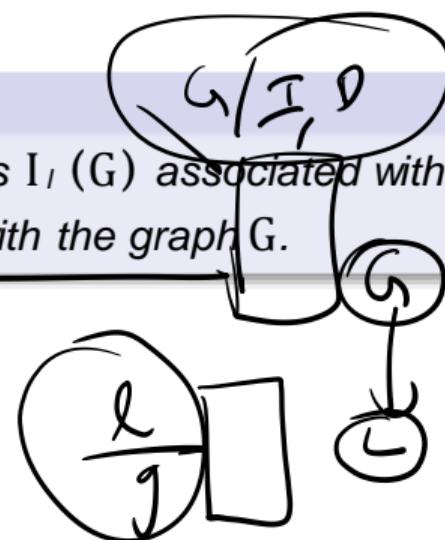
MARKOV BLANKET A node is conditionally independent of all other nodes in the Bayesian Network, given its parents, children and children's parents.

For each variable X_i : $(X_i \perp \text{other nodes} | \text{Pa}(X_i), \text{Ch}(X_i), \text{Pa}(\text{Ch}(X_i)))$

INDEPENDENCY MAP

THEOREM

A distribution P satisfies local independencies $I_1(G)$ associated with G if and only if P is representable as a set of CPDs associated with the graph G .



INDEPENDENCY MAP OR I-MAP

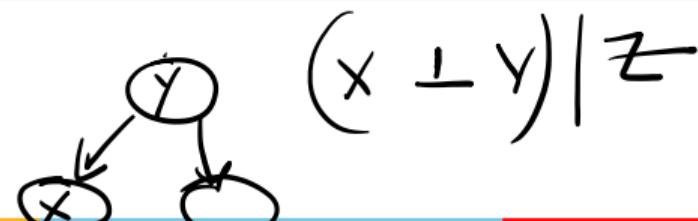
$$\cancel{P(X_1, X_2, \dots, X_n | Z)} \rightarrow P(X, Y | Z) = P(X | Z) P(Y | Z)$$



- P be a distribution over X .
- $I(P)$ be the set of independence assertions ($X \perp Y | Z$) that hold in P .
- Any independence that G asserts must also hold in P .

DEFINITION

G is called an **I-map** for P if $I_G \subseteq I(P)$.



$$P(X_0, Y_0) = P(X_0)P(Y_0) = 0.4 \times 0.2 = 0.08$$

I-MAP EXAMPLE 1

$$P(X_0) = 0.4$$

$$P(Y_0) = 0.2$$

X		Y		$P(X, Y)$
x^0	y^0	x^0	y^1	
x^0	y^0	0.08	0.32	$P(X_0, Y_0)$
x^1	y^0			
x^1	y^1	0.12	0.48	$P(X_1, Y_0)$

are X and Y
independent in P?

$$\underline{P(X, Y)}$$

$$P(X_0)P(Y_0) = P(X_0, Y_0)$$

$$P(X_0)P(Y_1) = P(X_0, Y_1)$$

$$P(X_1)P(Y_0) = P(X_1, Y_0)$$

$$X \perp Y$$

- Is G_ϕ : $X \perp Y$ an I-map of P ?

If we simply the graph



(no edge between X & Y)

I-MAP EXAMPLE 1

X	Y	$P(X, Y)$
x^0	y^0	0.08
x^0	y^1	0.32
x^1	y^0	0.12
x^1	y^1	0.48

- Is $G_\phi : X \perp Y$ an I-map of P ?

$\underbrace{G_\phi}_{\text{encodes the assumption}} \text{ that } X \perp Y.$

- $P(x^1) = 0.48 + 0.12 = 0.60$
- $P(y^1) = 0.32 + 0.48 = 0.80$
- $P(x^1, y^1) = 0.48 = P(x^1)P(y^1)$
- Hence X and Y are independent i.e $(X \perp Y)$
- $(X \perp Y) \in I(P)$.
- G_ϕ is an I-map of P .

Similarly $P(x^0, y^0) = P(x^0)P(y^0)$
 $P(x^0, y^1) = P(x^0)P(y^1)$ & $P(x^1, y^0) = P(x^1)P(y^0)$

I-MAP EXAMPLE 2

X	Y	$P(X, Y)$
x^0	y^0	0.4
x^0	y^1	0.3
x^1	y^0	0.2
x^1	y^1	0.1

$$P(X, Y) \stackrel{?}{=} P(X) P(Y)$$

- Is $G_\phi : X \perp\!\!\! \perp Y$ an I-map of P ?

I-MAP EXAMPLE 2

X	Y	$P(X, Y)$
x^0	y^0	0.4
x^0	y^1	0.3
x^1	y^0	0.2
x^1	y^1	0.1

- Is $G_\phi : X \perp\!\!\!\perp Y$ an I-map of P ?

- $P(x^1) = 0.2 + 0.1 = \cancel{0.3} \quad \textcircled{3}$
- $P(y^1) = \underline{0.3 + 0.1} = 0.4$
- $P(x^1, y^1) \neq P(x^1)P(y^1)$
- Hence X and Y are not independent.
- $(X \perp\!\!\!\perp Y) \notin I(P)$.
- G_ϕ is not an I-map of P .

STUDENT EXAMPLE - B^{Student}

We know independence assumptions in G

$$(D \perp I) \implies$$

$$\underline{P(D|I)} = P(D)$$

$$(L \perp I, D|G) \implies$$

$$\underline{P(L|I, D, G)} = P(L|G)$$

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

$$\underline{P(S|I, D, G, L)} = P(S|I)$$

$$P(I, D, G, S, L) = P(I)P(D|I)P(G|I, D)P(L|I, D, G)P(S|I, D, G, L)$$

(by chain rule)

$$= P(I)P(D)P(G|I, D)P(L|I, D, G)P(S|I, D, G, L)$$

$$= P(I)P(D)P(G|I, D)P(L|G)P(S|I, D, G, L)$$

$$= P(I)P(D)P(G|I, D)P(S|I)P(L|G)$$

$\prod_i P(x_i | Pa(x_i))$

A distribution P factorizes G if P can be represented as a chain rule of CPDs

P FACTORIZES OVER G

DEFINITION

Bayesian Network graph G over the variables X_1, X_2, \dots, X_n with a distribution P over the same space factorizes according to G if P can be expressed as a product of its CPDs.

$$P \text{ factorizes over } G \text{ if } P(X_1, X_2, \dots, X_n) = \prod_i^Y P(X_i | Pa^G(X_i)) \quad (3)$$

THEOREM

For a Bayesian Network graph G over the variables X_1, X_2, \dots, X_n with a distribution P over the same space and G is an I-map for P , then P factorizes G .

$\rightarrow G$ is an IMA for $P \Rightarrow P$ factorizes G

Proof of the theorem

Assume that X_1, X_2, \dots, X_n form a topological ordering (parents of X_i have subscripts smaller than i ; children of X_i have larger subscripts)

Using the chain rule

$$P(X_1, X_2, \dots, X_n) = P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots P(X_i | X_1, X_2, \dots, X_{i-1}) \dots$$

Proof of the theorem

Consider one of the factors, $l(x_i | x_1, x_2, \dots, x_{i-1})$

Since g is an I-map for P we have

$$(x_i \perp \text{NonDescendants } x_i) \mid P^{x_i}$$

Now all of x_i 's parents are in the set

$\{x_1, x_2, \dots, x_{i-1}\}$ and none of x_i 's descendants can be in this set

Proof of th. theorem

$\{x_1, x_2, \dots, x_{i-1}\} = Pa(x_i) \cup z$ where

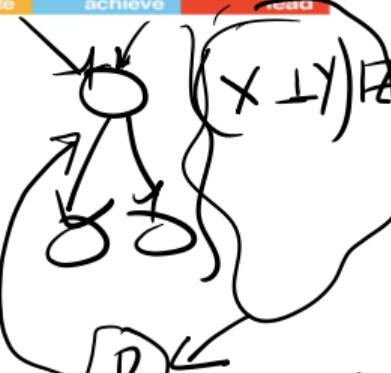
$z \subseteq \text{Non Descendants } x_i$

We know that $(x \perp z | Pa(x_i))$ which means that

$$P(x_i | Pa(x_i) \cup z) = P(x_i | Pa(x_i))$$

$$\prod P(x_i | Pa(x_i))$$

Now apply this to all the factors in the chain rule
to see that P factorizes G.



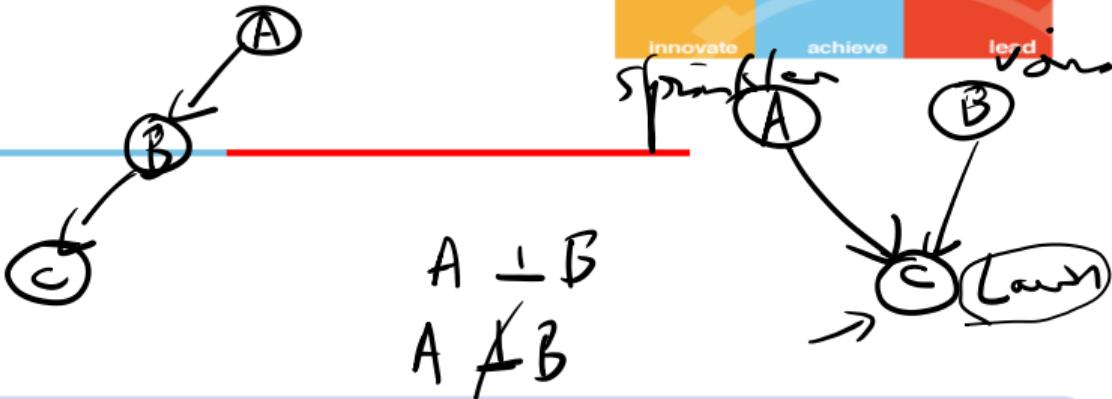
BAYESIAN NETWORK - ANOTHER DEFINITION

DEFINITION

A Bayesian network is a pair $B = (G, P)$ where P factorizes over G and where P is specified as a set of CPDs associated with G .

G IS AN I-MAP OF P

P(rain/sprinkler/lawn)



THEOREM

For a Bayesian Network graph G over the variables X_1, X_2, \dots, X_n with a distribution P over the same space and if P factorizes according to G, then G is an I-map for P.

G is an I-map for P \Rightarrow P factorizes according to G

P factorizes according to G \Rightarrow G is an IMAP for P

STUDENT EXAMPLE - B^{Student}

Given: P factors according to G

By chain rule $P(I, D, G, S, L) = \frac{P(I)P(D)P(G|I, D)P(S|I, D)P(L|G)}{P(I, D, G, S, L)}$

By definition $P(S|I, D, G, L) = \frac{P(S|I, D, G, L)}{P(I, D, G, L)}$

Marginalize over S $P(I, D, G, L) = \sum_S \underbrace{P(I, D, G, S, L)}_{S}$

$$= P(I)P(D)P(G|I, D)P(L|G)$$

$$\sum_S P(S|I)$$

$$P(A|I(B|A)) = P(A, B) = P(I)P(D)P(G|I, D)P(L|G)$$

$$\sum_S P(S|I)$$

STUDENT EXAMPLE - B^{Student}

$$\begin{aligned}
 P(S|I, D, G, L) &= \frac{P(I, D, G, S, L)}{P(I, D, G, L)} \\
 &= \frac{P(I)P(D)P(G|I, D)P(S|I)P(L|G)}{P(I)P(D)P(G|I, D)P(L|G)} \\
 &= P(S|I) \\
 &\quad (S \perp D, G, L | I)
 \end{aligned}$$

$$P(S|I, D, G, L) = P(S|I)$$

$$\boxed{S \perp D, G, L | I}$$

Independence assumption holds. G is an I-map for P.

REFERENCES

- 1 Probabilistic Graphical Models: Principles and Techniques by Daphne Koller and Nir Friedman. MIT Press. 2009
- 2 Artificial Intelligence: A Modern Approach (3rd Edition) by Stuart Russell, Peter Norvig
- 3 Mastering Probabilistic Graphical Models using Python by Ankur Ankan, Abhinash Panda. Packt Publishing 2015.
- 4 Learning in Graphical Models by Michael I. Jordan. MIT Press. 1999

Thank You!!!