

Readings:

K&F: 3.1, 3.2, 3.3

BN Semantics 1

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

September 15th, 2006

1

Let's start on BNs...

- Consider $P(X_i)$ *$k-1$ params.*
 - Assign probability to each $x_i \in \text{Val}(X_i)$
 - *k* Independent parameters
- Consider $P(X_1, \dots, X_n)$
 - How many independent parameters if $|\text{Val}(X_i)|=k$?
 $k^n - 1$ parameters

What if variables are independent?

■ What if variables are independent?

- $(X_i \perp X_j), \forall i, j$ ($\{X_1, X_3\} \perp \{X_7, X_9\}$)
- Not enough!!! (See homework 1 ☺) 1.7
- Must assume that $(\mathbf{X} \perp \mathbf{Y}), \forall \mathbf{X}, \mathbf{Y}$ subsets of $\{X_1, \dots, X_n\}$

■ Can write

□ $P(X_1, \dots, X_n) = \prod_{i=1 \dots n} P(X_i)$

■ How many independent parameters now?

$n \cdot (k-1)$ params

10-708 - ©Carlos Guestrin 2006

3

Conditional parameterization – two nodes

■ Grade is determined by Intelligence

$P(I)$ (I) $\text{val}[I] = \{H, VH\}$

$P(I, G) =$

$I \backslash G$	A	A+
H	0.01	0.15
VH	0.04	0.8

$P(G|I)$ (G) $\text{val}[G] = \{A+, A\}$

$P(I, G) = P(I) \cdot P(G|I)$

$P(I) =$

I	H	VH
	0.16	0.84

$P(G|I) =$

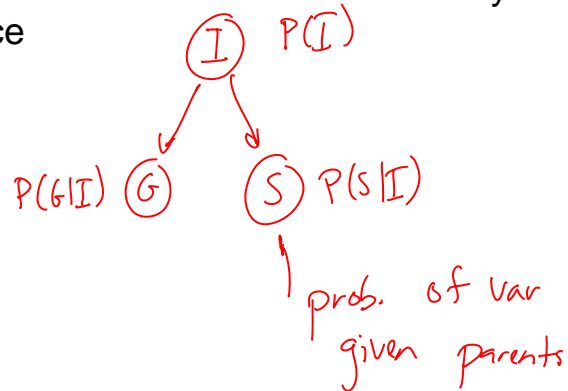
$I \backslash G$	A	A+
H	0.01 / 0.16	0.15 / 0.16
V	0.04 / 0.84	0.8 / 0.84

10-708 - ©Carlos Guestrin 2006

4

Conditional parameterization – three nodes

- Grade and SAT score are determined by Intelligence
- $(G \perp S \mid I)$

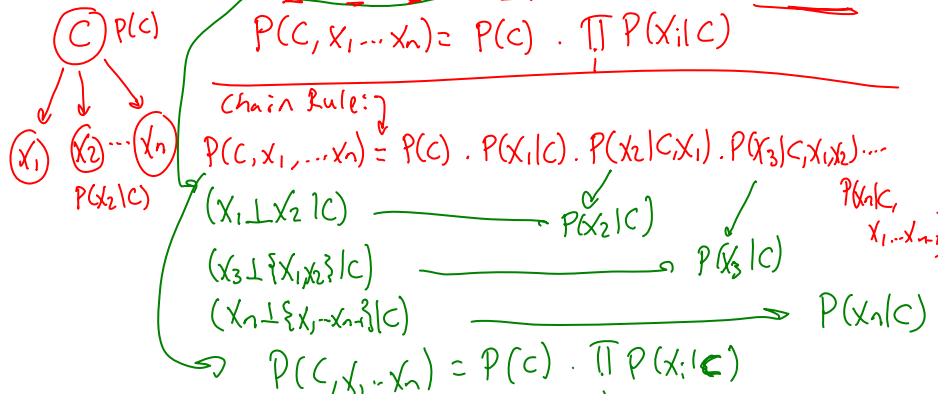


10-708 – ©Carlos Guestrin 2006

5

The naïve Bayes model – Your first real Bayes Net

- Class variable: C
- Evidence variables: X_1, \dots, X_n
- assume that $(X \perp Y \mid C), \forall X, Y$ subsets of $\{X_1, \dots, X_n\}$



10-708 – ©Carlos Guestrin 2006

6

What you need to know (From last class)

- Basic definitions of probabilities
- Independence
- Conditional independence
- The chain rule
- Bayes rule
- Naïve Bayes

10-708 – ©Carlos Guestrin 2006

7

Announcements

- Homework 1:
 - Out yesterday
 - Due September 27th – **beginning of class!**
 - It's hard – start early, ask questions
- Collaboration policy
 - OK to discuss in groups
 - Tell us on your paper who you talked with
 - Each person must write their **own unique paper**
 - No searching the web, papers, etc. for answers, we trust you want to learn
- Upcoming recitation
 - Monday 5:30-7pm in Wean 4615A – Matlab Tutorial
- Don't forget to register to the mailing list at:
 - <https://mailman.srv.cs.cmu.edu/mailman/listinfo/10708-announce>

10-708 – ©Carlos Guestrin 2006

8

This class

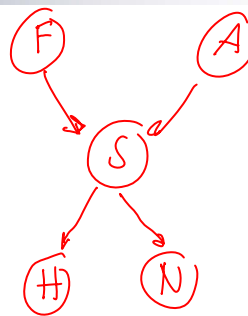
- We've heard of Bayes nets, we've played with Bayes nets, we've even used them in your research
- This class, we'll learn the semantics of BNs, relate them to independence assumptions encoded by the graph

10-708 - ©Carlos Guestrin 2006

9

Causal structure

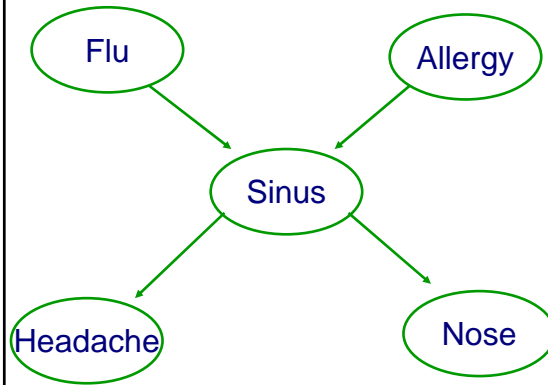
- Suppose we know the following:
 - The flu causes sinus inflammation
 - Allergies cause sinus inflammation
 - Sinus inflammation causes a runny nose
 - Sinus inflammation causes headaches
- How are these connected?



10-708 - ©Carlos Guestrin 2006

10

Possible queries



- Inference

$$P(F|H=t)$$

- Most probable explanation

$$\max_{f,a,s} P(F=f, A=a, S=s|H=t, N=f)$$

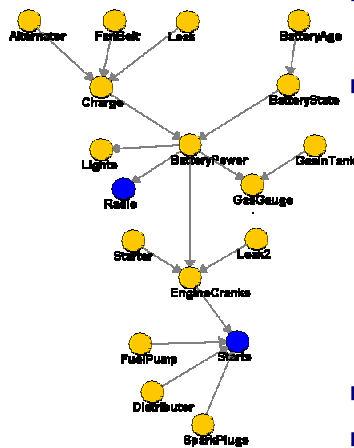
- Active data collection

I say $H=t$
tell me the value of N ?

10-708 – ©Carlos Guestrin 2006

11

Car starts BN



- 18 binary attributes

- Inference

$$P(\text{BatteryAge} | \text{Starts}=f)$$

- 2^{18} terms, why so fast?

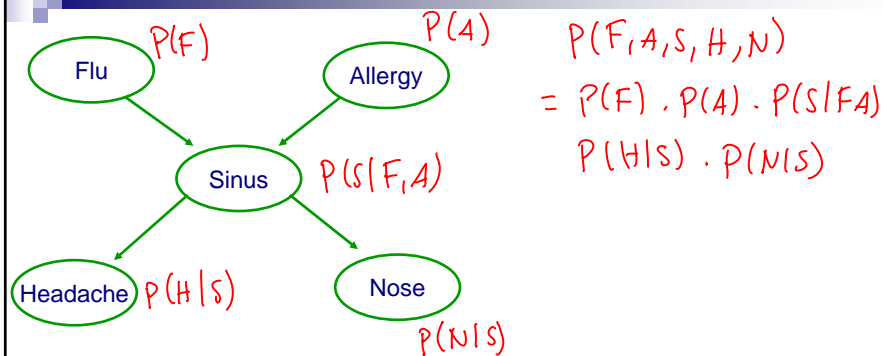
- Not impressed?

$$P(\text{HailFinder BN}) - \text{more than } 3^{54} =$$

58149737003040059690390169 terms

10-708 – ©Carlos Guestrin 2006

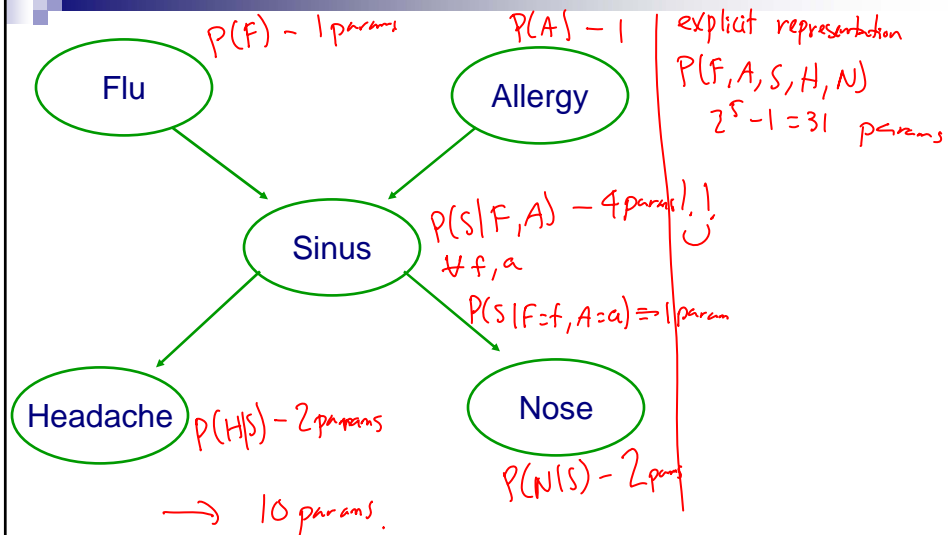
Factored joint distribution - Preview



10-708 - ©Carlos Guestrin 2006

13

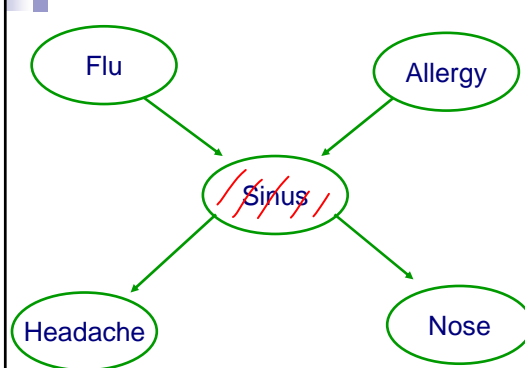
Number of parameters (binary vars)



10-708 - ©Carlos Guestrin 2006

14

Key: Independence assumptions



$(F \perp A)$
 not $(F \perp N)$
 $(F \perp N | S)$
 not $(H \perp N)$
 $(H \perp N | S)$

Knowing sinus separates the variables from each other

10-708 - ©Carlos Guestrin 2006

15

(Marginal) Independence

- Flu and Allergy are (marginally) independent

$$F \perp A \Rightarrow P(F, A) = P(F) \cdot P(A)$$

- More Generally:

Flu = t	0.1
Flu = f	0.9

Allergy = t	0.2
Allergy = f	0.8

	Flu = t	Flu = f
Allergy = t	0.1 × 0.2	0.2 × 0.9
Allergy = f		

10-708 - ©Carlos Guestrin 2006

16

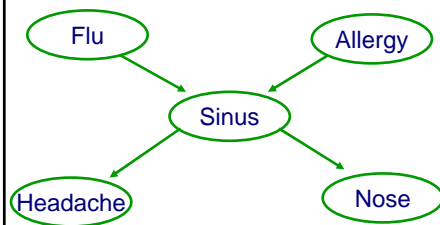
Conditional independence

- Flu and Headache are not (marginally) independent ^{not $(F \perp H)$}
- Flu and Headache are independent given Sinus infection $(F \perp H | S) \rightarrow P(F|S, H) = P(F|S)$
- More Generally: $(X \perp Y | Z) \Rightarrow P(X|Y, Z) = P(X|Z)$

10-708 - ©Carlos Guestrin 2006

17

The independence assumption



$$(F \perp A | \emptyset) \equiv (F \perp A)$$

$$(N \perp \{F, A, H\} | S)$$

Local Markov Assumption:

A variable X is independent of its non-descendants given its parents and only its parents $(X_i \perp \text{NonDescendants}_{X_i} | \text{Pa}_{X_i})$

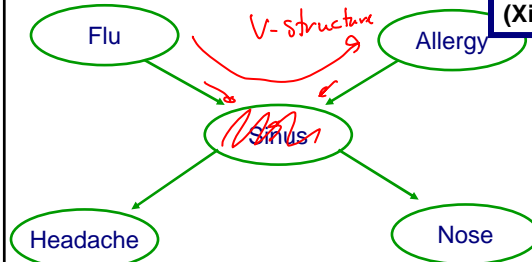
10-708 - ©Carlos Guestrin 2006

18

Explaining away

Local Markov Assumption:

A variable X is independent of its non-descendants given its parents *and only its parents*
 $(X_i \perp \text{NonDescendants}_{X_i} \mid \text{Pa}_{X_i})$



$$(F \perp A)$$

$$P(F=t) \leq P(F=t \mid S=t)$$

$$P(A=t) \leq P(A=t \mid S=t)$$

$$P(A=t \mid S=t, F=t) \leq$$

$$P(A=t) \leq P(A=t \mid S=t)$$

$$(F \perp A)$$

$$\text{not } (F \perp A \mid S)$$

$$\text{not } (F \perp A \mid H)$$

S was the center of a V-structure
 H was the descendent of " "

10-708 - ©Carlos Guestrin 2006

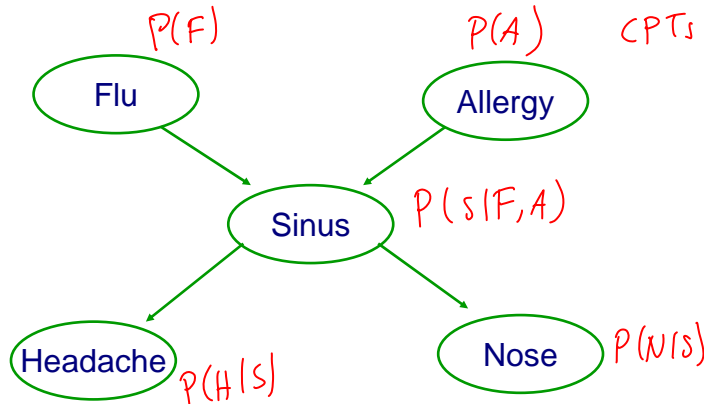
19

What about probabilities?

discrete setting

Conditional probability tables (CPTs)

Also called "Factors": in a BN a CPT is a factor

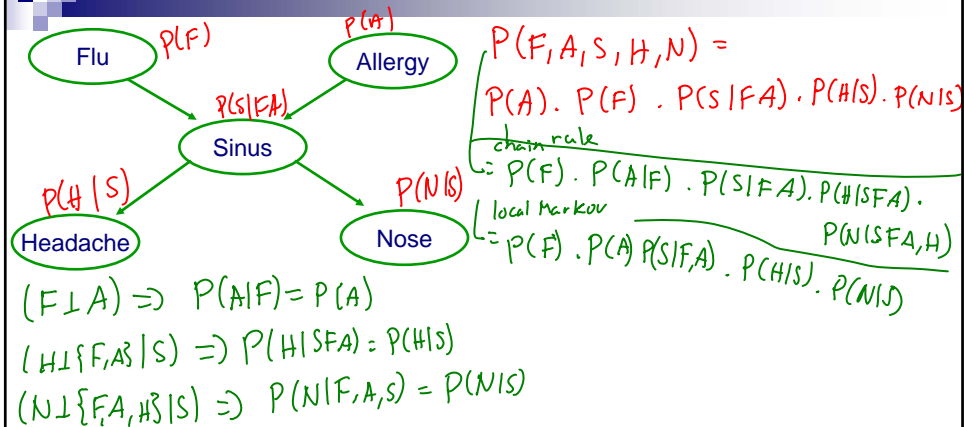


10-708 - ©Carlos Guestrin 2006

20

Joint distribution

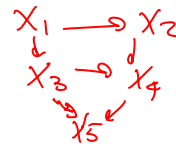
→ Chain rule
→ Local Markov Assumption



Why can we decompose? ^{local} Markov Assumption!

A general Bayes net

- Set of random variables $\{X_1, \dots, X_n\}$
- Directed acyclic graph ^{DAG}
- CPTs $P(X_i | Pa_{X_i})$; $Pa_{X_i} \subseteq \{X_1, \dots, X_n\}$



- Joint distribution:

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

- Local Markov Assumption:

- A variable X is independent of its non-descendants given its parents – $(X \perp NonDescendantsX | PaX)$
^{randomly its parents}

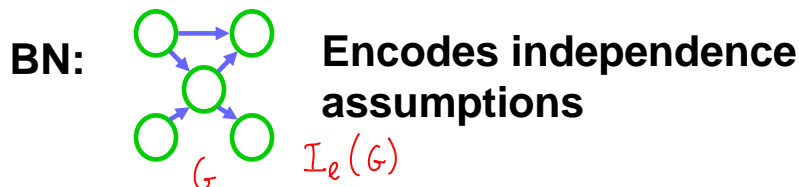
Questions????

- What distributions can be represented by a BN?
- What BNs can represent a distribution?
- What are the independence assumptions encoded in a BN?
 - in addition to the local Markov assumption

10-708 – ©Carlos Guestrin 2006

23

Today: The Representation Theorem – Joint Distribution to BN



If conditional independencies in BN are subset of conditional independencies in P

Obtain

Joint probability distribution:

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

$I_e(G) \subseteq I(P)$ ✓ indep. in P

10-708 – ©Carlos Guestrin 2006

24

Today: The Representation Theorem – BN to Joint Distribution

BN:  **Encodes independence assumptions**

If joint probability distribution:

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

Obtain →

Then conditional independencies in BN are subset of conditional independencies in P

$$I_e(G) \subseteq I(P)$$

10-708 – ©Carlos Guestrin 2006

25

Let's start proving it for naïve Bayes – From joint distribution to BN

- Independence assumptions:
 - X_i independent given C
- Let's assume that P satisfies independencies must prove that P factorizes according to BN:
 - $P(C, X_1, \dots, X_n) = P(C) \prod_i P(X_i | C)$
- Use chain rule!

$$\begin{aligned} P(C, X_1, \dots, X_n) &= P(C) \cdot P(X_1 | C) \cdot P(X_2 | X_1, C) \cdots \\ &= P(C) \cdot P(X_1 | C) \cdot P(X_2 | C) \cdots \end{aligned}$$

10-708 – ©Carlos Guestrin 2006

26

Let's start proving it for naïve Bayes – From BN to joint distribution 1

- Let's assume that P factorizes according to the BN:

$$\square P(C, X_1, \dots, X_n) = P(C) \prod_i P(X_i | C)$$

- Prove the independence assumptions:

$$\square X_i \text{ independent given } C$$

$$\square \text{Actually, } (X \perp Y | C), \forall X, Y \text{ subsets of } \{X_1, \dots, X_n\}$$

$$P(X|Y, C) = P(X|C) \quad \text{equiv.} \quad P(X, Y|C) = P(X|C) \cdot P(Y|C)$$

$$X = \{X_1, X_2\} \quad Y = \{X_3, X_4\} \quad P(X_1, \dots, X_4 | C) = P(X_1, X_2 | C) \cdot P(X_3, X_4 | C)$$

$$P(X_1, \dots, X_4 | C) = \frac{P(C)}{P(C)} \cdot P(X_1 | C) \dots P(X_4 | C) = \underbrace{P(X_1 | C) \cdot P(X_2 | C)}_{P(X_1, X_2 | C)} \cdot \underbrace{P(X_3 | C) \cdot P(X_4 | C)}_{P(X_3, X_4 | C)}$$

$$P(X_1, X_2 | C) = \sum_{x_3, x_4} P(X_1, \dots, X_4 | C)$$

$$\{X_1, X_2\} \perp \{X_3, X_4\} | C \quad \text{g.e.d.}$$

10-708 – ©Carlos Guestrin 2006

27

Let's start proving it for naïve Bayes – From BN to joint distribution 2

- Let's consider a simpler case

$$\square \text{Grade and SAT score are determined by Intelligence}$$

$$\square P(I, G, S) = P(I)P(G|I)P(S|I)$$

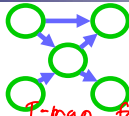
$$\square \text{Prove that } P(G, S|I) = P(G|I) P(S|I)$$

10-708 – ©Carlos Guestrin 2006

28

Today: The Representation Theorem

BN:



Encodes independence assumptions

I-map: G is an I-map for P

If conditional independencies in BN are subset of conditional independencies in P
 $I_c(G) \subseteq I(P)$

Obtain

Joint probability distribution:

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

If joint probability distribution:

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

Obtain

Then conditional independencies in BN are subset of conditional independencies in P

10-708 - ©Carlos Guestrin 2006

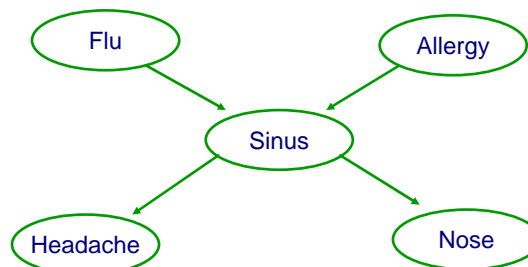
29

Local Markov assumption & I-maps

- Assumptions
- Local independence assumptions in BN structure G : *$I_c(G)$*

- Truth:
- Independence assertions of P : *$I(P)$*

- BN structure G is an **I-map** (independence map) if: *$I_c(G) \subseteq I(P)$*



Local Markov Assumption:

A variable X is independent of its non-descendants given its parents
 $(X_i \perp \text{NonDescendants}_{X_i} | \text{Pa}_{X_i})$

10-708 - ©Carlos Guestrin 2006

30

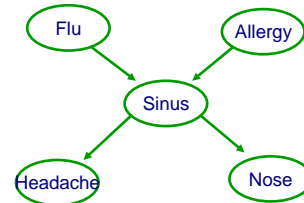
Factorized distributions

■ Given

- Random vars X_1, \dots, X_n
- P distribution over vars
- BN structure G over same vars

■ P factorizes according to G if

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



10-708 – ©Carlos Guestrin 2006

31

BN Representation Theorem – I-map to factorization

If conditional independencies in BN are subset of conditional independencies in P

Obtain

Joint probability distribution:

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

G is an I-map of P

P factorizes according to G

10-708 – ©Carlos Guestrin 2006

32

BN Representation Theorem – I-map to factorization: **Proof**

**G is an
I-map of P**

Obtain

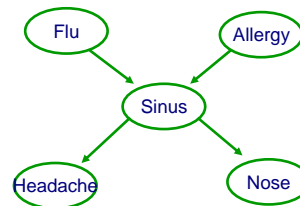
**P factorizes
according to G**

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

ALL YOU NEED:

Local Markov Assumption:

A variable X is independent
of its non-descendants given its parents
($X_i \perp \text{NonDescendants}_{X_i} \mid \text{Pa}_{X_i}$)



10-708 – ©Carlos Guestrin 2006

33

Defining a BN

- Given a set of variables and conditional independence assertions of P
- Choose an ordering on variables, e.g., X_1, \dots, X_n
- For $i = 1$ to n
 - Add X_i to the network
 - Define parents of X_i , Pa_{X_i} , in graph as the minimal subset of $\{X_1, \dots, X_{i-1}\}$ such that local Markov assumption holds – X_i independent of rest of $\{X_1, \dots, X_{i-1}\}$, given parents Pa_{X_i}
 - Define/learn CPT – $P(X_i \mid \text{Pa}_{X_i})$

10-708 – ©Carlos Guestrin 2006

34

BN Representation Theorem – Factorization to I-map

If joint probability
distribution:

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

Obtain

Then conditional
independencies
in BN are subset of
conditional
independencies in P

**P factorizes
according to G**

G is an I-map of P

10-708 – ©Carlos Guestrin 2006

35

BN Representation Theorem – Factorization to I-map: **Proof**

If joint probability
distribution:

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

Obtain

Then conditional
independencies
in BN are subset of
conditional
independencies in P

**P factorizes
according to G**

G is an I-map of P

Homework 1!!!! 😊

10-708 – ©Carlos Guestrin 2006

36

The BN Representation Theorem

If conditional independencies in BN are subset of conditional independencies in P

Obtain

Joint probability distribution:

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

Important because:
Every P has at least one BN structure G

If joint probability distribution:

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

Obtain

Then conditional independencies in BN are subset of conditional independencies in P

Important because:
Read independencies of P from BN structure G

10-708 – ©Carlos Guestrin 2006

37

Independencies encoded in BN

- We said: All you need is the local Markov assumption
 - $(X_i \perp \text{NonDescendants}_{X_i} | \text{Pa}_{X_i})$
- But then we talked about other (in)dependencies
 - e.g., explaining away
- What are the independencies encoded by a BN?
 - Only assumption is local Markov
 - But many others can be derived using the algebra of conditional independencies!!!

10-708 – ©Carlos Guestrin 2006

38

Understanding independencies in BNs

– BNs with 3 nodes

Local Markov Assumption:

A variable X is independent of its non-descendants given its parents

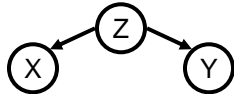
Indirect causal effect:



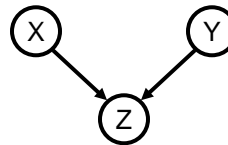
Indirect evidential effect:



Common cause:



Common effect:

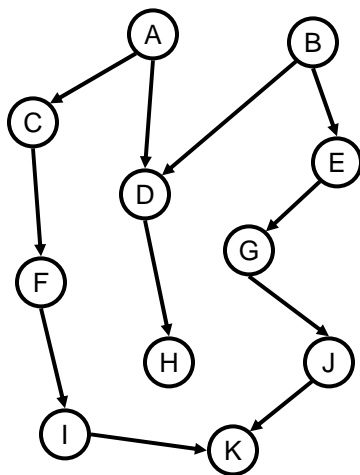


10-708 – ©Carlos Guestrin 2006

39

Understanding independencies in BNs

– Some examples

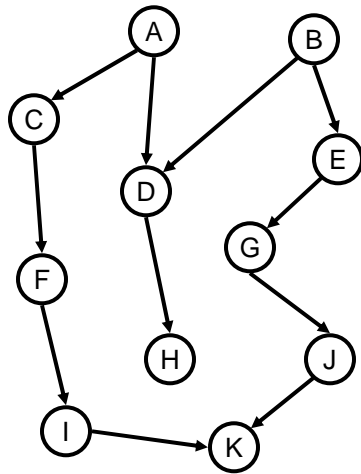


10-708 – ©Carlos Guestrin 2006

40

Understanding independencies in BNs

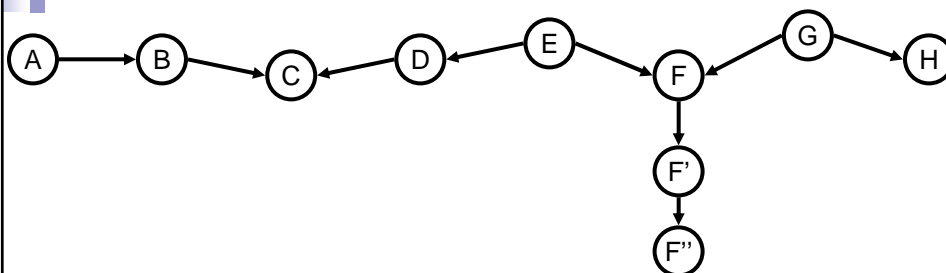
– Some more examples



10-708 – ©Carlos Guestrin 2006

41

An active trail – Example



When are A and H independent?

10-708 – ©Carlos Guestrin 2006

42

Active trails formalized

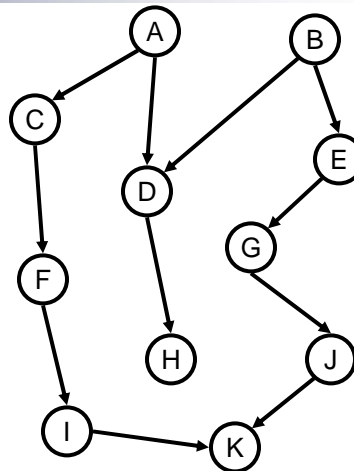
- A trail $X_1 - X_2 - \dots - X_k$ is an **active trail** when variables $\mathbf{O} \subseteq \{X_1, \dots, X_n\}$ are observed if for each consecutive triplet in the trail:
 - $X_{i-1} \rightarrow X_i \rightarrow X_{i+1}$, and X_i is **not observed** ($X_i \notin \mathbf{O}$)
 - $X_{i-1} \leftarrow X_i \leftarrow X_{i+1}$, and X_i is **not observed** ($X_i \notin \mathbf{O}$)
 - $X_{i-1} \leftarrow X_i \rightarrow X_{i+1}$, and X_i is **not observed** ($X_i \notin \mathbf{O}$)
 - $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$, and X_i is **observed** ($X_i \in \mathbf{O}$), or **one of its descendants**

10-708 – ©Carlos Guestrin 2006

43

Active trails and independence?

- **Theorem:** Variables X_i and X_j are independent given $\mathbf{Z} \subseteq \{X_1, \dots, X_n\}$ if there is **no active trail** between X_i and X_j when variables $\mathbf{Z} \subseteq \{X_1, \dots, X_n\}$ are observed



10-708 – ©Carlos Guestrin 2006

44

More generally:

Soundness of d-separation

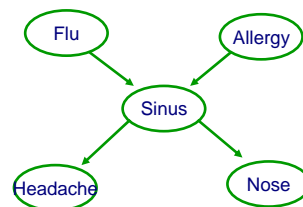
- Given BN structure G
- Set of independence assertions obtained by d-separation:
 - $I(G) = \{(X \perp Y | Z) : \text{d-sep}_G(X; Y | Z)\}$
- **Theorem: Soundness of d-separation**
 - If P factorizes over G then $I(G) \subseteq I(P)$
- **Interpretation:** d-separation only captures true independencies
- Proof discussed when we talk about undirected models

10-708 – ©Carlos Guestrin 2006

45

Adding edges doesn't hurt

- **Theorem:** Let G be an I-map for P , any DAG G' that includes the same directed edges as G is also an I-map for P .
- **Proof sketch:**



10-708 – ©Carlos Guestrin 2006

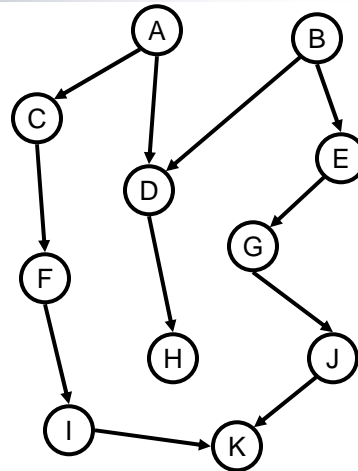
46

Existence of dependency when not d-separated

- **Theorem:** If X and Y are not d-separated given Z , then X and Y are dependent given Z under some P that factorizes over G

- **Proof sketch:**

- Choose an active trail between X and Y given Z
- Make this trail dependent
- Make all else uniform (independent) to avoid “canceling” out influence



10-708 – ©Carlos Guestrin 2006

47

More generally: Completeness of d-separation

- **Theorem: Completeness of d-separation**

- For “almost all” distributions that P factorize over to G , we have that $I(G) = I(P)$
- “almost all” distributions: except for a set of measure zero of parameterizations of the CPTs (assuming no finite set of parameterizations has positive measure)

- **Proof sketch:**

10-708 – ©Carlos Guestrin 2006

48

Interpretation of completeness

■ Theorem: Completeness of d-separation

- For “almost all” distributions that P factorize over to G , we have that $I(G) = I(P)$
- BN graph is usually sufficient to capture all independence properties of the distribution!!!!
- But only for complete independence:
 - $P \models (X=x \perp Y=y \mid Z=z), \forall x \in \text{Val}(X), y \in \text{Val}(Y), z \in \text{Val}(Z)$
- Often we have context-specific independence (CSI)
 - $\exists x \in \text{Val}(X), y \in \text{Val}(Y), z \in \text{Val}(Z): P \models (X=x \perp Y=y \mid Z=z)$
 - Many factors may affect your grade
 - But if you are a frequentist, all other factors are irrelevant ☺

10-708 – ©Carlos Guestrin 2006

49

What you need to know

- Independence & conditional independence
- Definition of a BN
- The representation theorems
 - Statement
 - Interpretation
- d-separation and independence
 - soundness
 - existence
 - completeness

10-708 – ©Carlos Guestrin 2006

50

Acknowledgements

- JavaBayes applet

- <http://www.pmr.poli.usp.br/ltd/Software/javabayes/Home/index.html>