# VC Dimension

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

October 29th, 2007

1

---

# What about continuous hypothesis spaces?

$$\text{error}_{true}(h) \leq \text{error}_{train}(h) + \sqrt{\frac{\ln|H| + \ln\frac{1}{\delta}}{2m}}$$
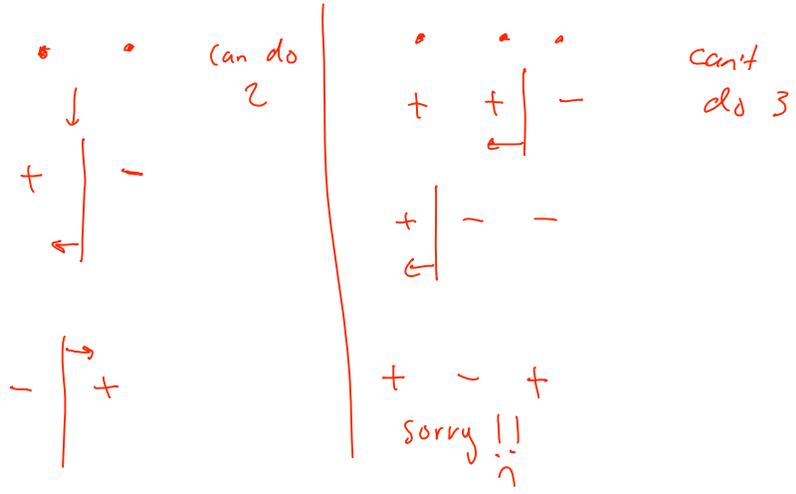
- Continuous hypothesis space:
  - $|H| = \infty$
  - Infinite variance???

- **As with decision trees, only care about the maximum number of points that can be classified exactly!**

2

1

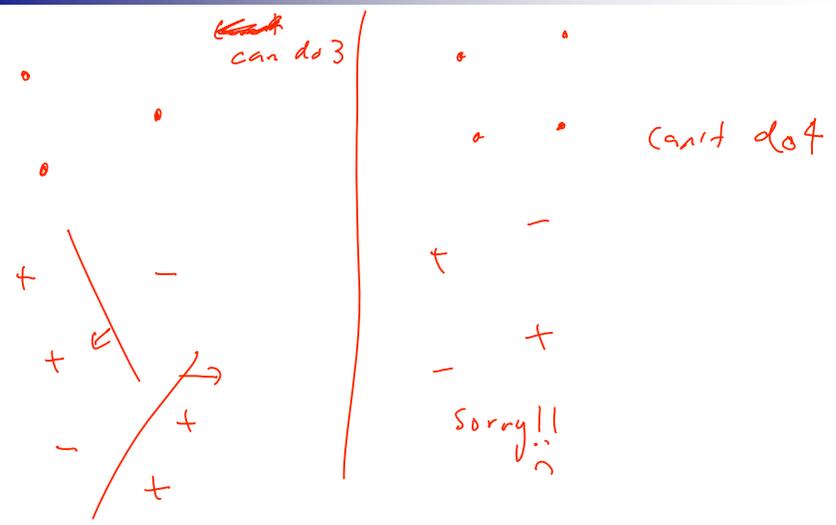# How many points can a linear boundary classify exactly? (1-D)

3

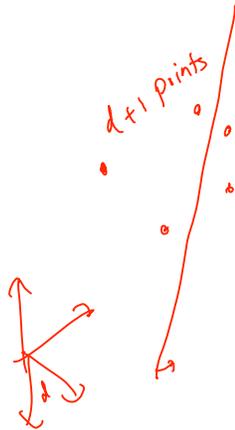# How many points can a linear boundary classify exactly? (2-D)

4

2

# How many points can a linear boundary classify exactly? (d-D)

d+1 points

can do d+1 points

how many parameters in a linear classifier in d-dimensions?

$$w_0 + \sum_{i=1}^{d} w_i x_i$$

d+1

# PAC bound using VC dimension

e.g., linear classifiers

- **Number of training points that can be classified exactly is VC dimension!!!**
  - □ **Measures relevant size of hypothesis space, as with decision trees with k leaves**

may be continuous

$$\text{error}_{true}(h) \leq \text{error}_{train}(h) + \sqrt{\frac{VC(H)\left(\ln\frac{2m}{VC(H)} + 1\right) + \ln\frac{4}{\delta}}{m}}$$

only depend on VC(H) not on |H|

# Shattering a set of points

*Definition:* a **dichotomy** of a set $S$ is a partition of $S$ into two disjoint subsets.

*Definition:* a set of instances $S$ is **shattered** by hypothesis space $H$ if and only if for every dichotomy of $S$ there exists some hypothesis in $H$ consistent with this dichotomy.

7

# VC dimension

*Definition:* The **Vapnik-Chervonenkis dimension**, $VC(H)$, of hypothesis space $H$ defined over instance space $X$ is the size of the largest finite subset of $X$ shattered by $H$. If arbitrarily large finite sets of $X$ can be shattered by $H$, then $VC(H) \equiv \infty$.

8

4

# PAC bound using VC dimension

- **Number of training points that can be classified exactly is VC dimension!!!**
  - ☐ **Measures relevant size of hypothesis space, as with decision trees with k leaves**
  - ☐ **Bound for infinite dimension hypothesis spaces:**

$$\text{error}_{true}(h) \leq \text{error}_{train}(h) + \sqrt{\frac{VC(H)\left(\ln\frac{2m}{VC(H)} + 1\right) + \ln\frac{4}{\delta}}{m}}$$

9

# Examples of VC dimension

$$\text{error}_{true}(h) \leq \text{error}_{train}(h) + \sqrt{\frac{VC(H)\left(\ln\frac{2m}{VC(H)} + 1\right) + \ln\frac{4}{\delta}}{m}}$$

- Linear classifiers:
  - ☐ VC(H) = d+1, for *d* features plus constant term *b*

- Neural networks
  - ☐ VC(H) = #parameters
  - ☐ Local minima means NNs will probably not find best parameters

- 1-Nearest neighbor?

10

5

# Another VC dim. example - What can we shatter?

- What's the VC dim. of decision stumps in 2d?

**11**

# Another VC dim. example - What can't we shatter?

- What's the VC dim. of decision stumps in 2d?

**12**

# What you need to know

- Finite hypothesis space
  - □ Derive results
  - □ Counting number of hypothesis
  - □ Mistakes on Training data
- Complexity of the classifier depends on number of points that can be classified exactly
  - □ Finite case – decision trees
  - □ Infinite case – VC dimension
- Bias-Variance tradeoff in learning theory
- Remember: will your algorithm find best classifier?

**13**

# Bayesian Networks – Representation
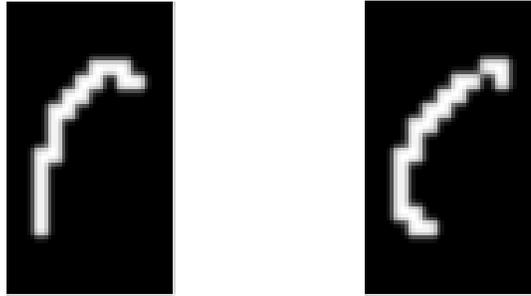
Machine Learning – 10701/15781

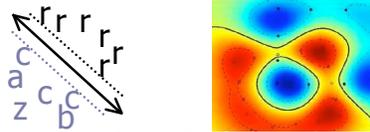Carlos Guestrin

Carnegie Mellon University

October 29th, 2007

**14**

# Handwriting recognition



Character recognition, e.g., kernel SVMs

15

# Webpage classification



→ Company home page
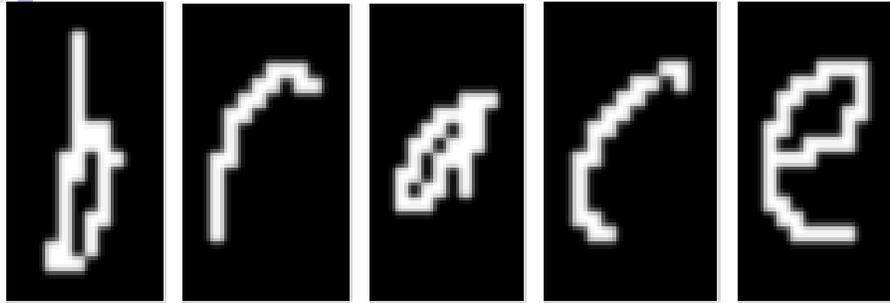
  vs

Personal home page

  vs

University home page

  vs

…

16

# Handwriting recognition 2

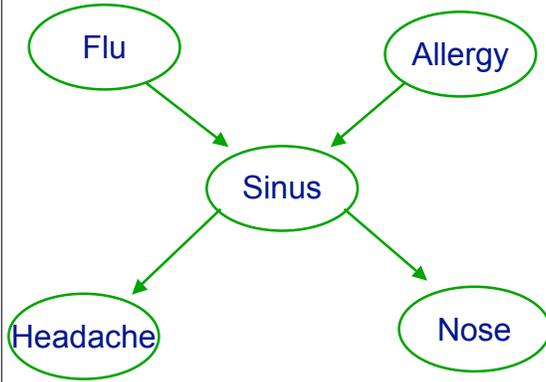17

# Webpage classification 2

18

9

# Today – Bayesian networks

- One of the most exciting advancements in statistical AI in the last 10-15 years
- Generalizes naïve Bayes and logistic regression classifiers
- Compact representation for exponentially-large probability distributions
- Exploit conditional independencies

19

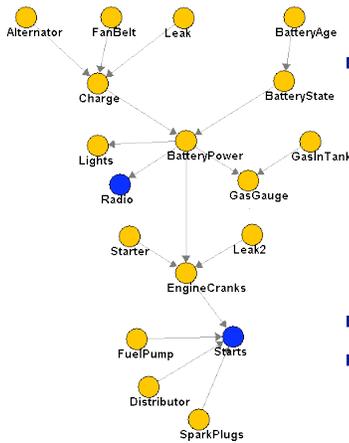# Causal structure

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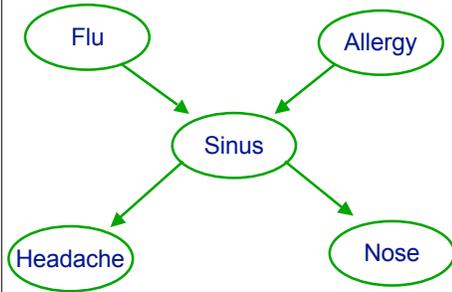- How are these connected?

20

# Possible queries

Flu

Allergy

Sinus

Headache

Nose

- Inference

- Most probable explanation

- Active data collection

# Car starts BN

Alternator   FanBelt   Leak        BatteryAge

Charge                 BatteryState

Lights   BatteryPower   GasInTank

Radio        GasGauge

Starter    Leak2

EngineCranks

FuelPump   Starts

Distributor

SparkPlugs

- 18 binary attributes

- Inference
  - $P(BatteryAge|Starts=f)$

- $2^{16}$ terms, why so fast?
- Not impressed?
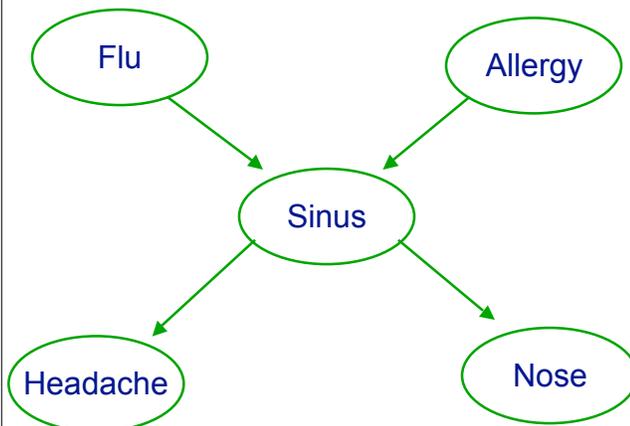  - HailFinder BN – more than $3^{54}$ = 58149737003040059690390169 terms
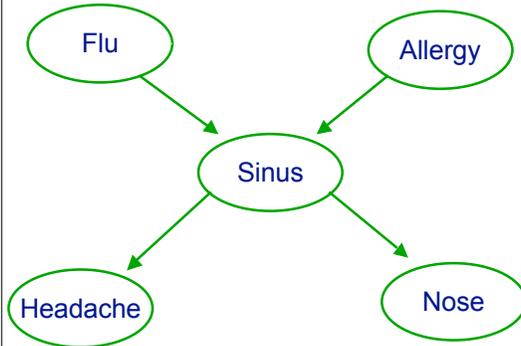
11

# Factored joint distribution - Preview

Flu → Sinus ← Allergy

Sinus → Headache

Sinus → Nose

23

# Number of parameters

Flu → Sinus ← Allergy

Sinus → Headache

Sinus → Nose

24

# Key: Independence assumptions



Flu     Allergy

Sinus

Headache     Nose

Knowing sinus separates the variables from each other

# (Marginal) Independence

- Flu and Allergy are (marginally) independent

| Flu = t | |
|---|---|
| Flu = f | |

- More Generally:

| Allergy = t | |
|---|---|
| Allergy = f | |

| | Flu = t | Flu = f |
|---|---|---|
| Allergy = t | | |
| Allergy = f | | |

# Marginally independent random variables

- **Sets** of variables **X**, **Y**
- X is independent of Y if
  - $P \vDash (\mathbf{X}{=}\mathbf{x} \perp \mathbf{Y}{=}\mathbf{y})$, $\forall$ $\mathbf{x} \in \mathrm{Val}(\mathbf{X})$, $\mathbf{y} \in \mathrm{Val}(\mathbf{Y})$

- Shorthand:
  - **Marginal independence:** $P \vDash (\mathbf{X} \perp \mathbf{Y})$

- **Proposition:** $P$ statisfies $(\mathbf{X} \perp \mathbf{Y})$ if and only if
  - $P(\mathbf{X},\mathbf{Y}) = P(\mathbf{X})\, P(\mathbf{Y})$

©2005-2007 Carlos Guestrin

27

# Conditional independence

- Flu and Headache are not (marginally) independent

- Flu and Headache are independent given Sinus infection

- More Generally:

©2005-2007 Carlos Guestrin

28

# Conditionally independent random variables

- **Sets** of variables **X**, **Y**, **Z**
- X is independent of Y given Z if
  - $P \vDash (\mathbf{X}{=}\mathbf{x} \perp \mathbf{Y}{=}\mathbf{y} | \mathbf{Z}{=}\mathbf{z})$, $\forall$ $\mathbf{x} \in$ Val(**X**), $\mathbf{y} \in$ Val(**Y**), $\mathbf{z} \in$ Val(**Z**)

- Shorthand:
  - **Conditional independence:** $P \vDash (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$
  - For $P \vDash (\mathbf{X} \perp \mathbf{Y} | \emptyset)$, write $P \vDash (\mathbf{X} \perp \mathbf{Y})$

- **Proposition:** *P* statisfies $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$ if and only if
  - P(**X**,**Y**|**Z**) = P(**X**|**Z**) P(**Y**|**Z**)

# Properties of independence

- **Symmetry:**
  - $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z}) \Rightarrow (\mathbf{Y} \perp \mathbf{X} | \mathbf{Z})$
- **Decomposition:**
  - $(\mathbf{X} \perp \mathbf{Y},\mathbf{W} | \mathbf{Z}) \Rightarrow (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$
- **Weak union:**
  - $(\mathbf{X} \perp \mathbf{Y},\mathbf{W} | \mathbf{Z}) \Rightarrow (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z},\mathbf{W})$
- **Contraction:**
  - $(\mathbf{X} \perp \mathbf{W} | \mathbf{Y},\mathbf{Z})$ & $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z}) \Rightarrow (\mathbf{X} \perp \mathbf{Y},\mathbf{W} | \mathbf{Z})$
- **Intersection:**
  - $(\mathbf{X} \perp \mathbf{Y} | \mathbf{W},\mathbf{Z})$ & $(\mathbf{X} \perp \mathbf{W} | \mathbf{Y},\mathbf{Z}) \Rightarrow (\mathbf{X} \perp \mathbf{Y},\mathbf{W} | \mathbf{Z})$
  - Only for positive distributions!
  - $P(\alpha)>0$, $\forall \alpha$, $\alpha \neq \emptyset$

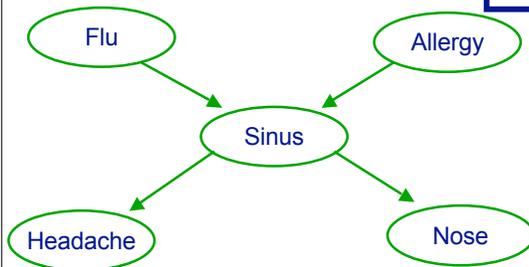# **The** independence assumption

Flu → Sinus ← Allergy
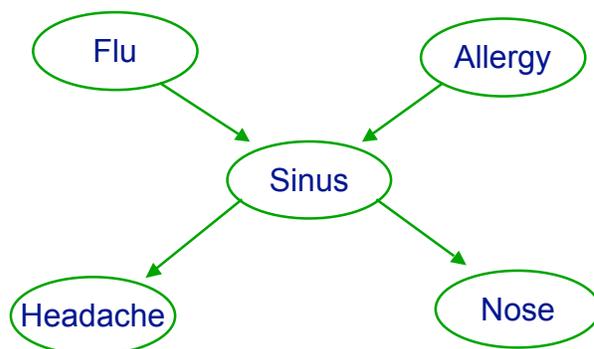Sinus → Headache
Sinus → Nose

**Local Markov Assumption:**
A variable X is independent of its non-descendants given its parents

---

# Explaining away

**Local Markov Assumption:**
A variable X is independent of its non-descendants given its parents

Flu → Sinus ← Allergy
Sinus → Headache
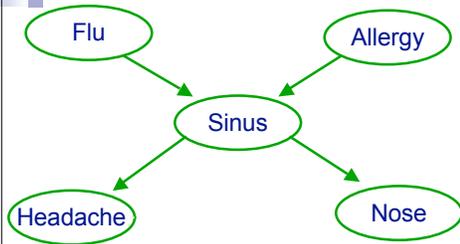Sinus → Nose

# Naïve Bayes revisited

> **Local Markov Assumption:**
> A variable X is independent of its non-descendants given its parents

33

# What about probabilities?
# Conditional probability tables (CPTs)

Flu → Sinus ← Allergy

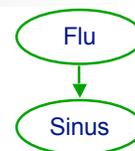Sinus → Headache

Sinus → Nose

34

17

# Joint distribution



**Why can we decompose? Markov Assumption!**

---

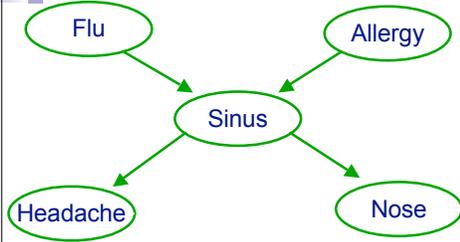# The chain rule of probabilities

- P(A,B) = P(A)P(B|A)



- More generally:
  - $P(X_1,\ldots,X_n) = P(X_1) \cdot P(X_2|X_1) \cdot \ldots \cdot P(X_n|X_1,\ldots,X_{n-1})$

36

18

# Chain rule & Joint distribution

**Local Markov Assumption:**
A variable X is independent of its non-descendants given its parents

Flu → Sinus
Allergy → Sinus
Sinus → Headache
Sinus → Nose

37

# Two (trivial) special cases

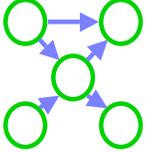**Edgeless graph**

**Fully-connected graph**

38

## The Representation Theorem – Joint Distribution to BN

**BN:** Encodes independence assumptions

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

Obtain →

**Joint probability distribution:**

$$P(X_1, \ldots, X_n) = \prod_{i=1}^{n} P\left(X_i \mid \mathbf{Pa}_{X_i}\right)$$

39

---

## Real Bayesian networks applications

- Diagnosis of lymph node disease
- Speech recognition
- Microsoft office and Windows
  - http://www.research.microsoft.com/research/dtg/
- Study Human genome
- Robot mapping
- Robots to identify meteorites to study
- Modeling fMRI data
- Anomaly detection
- Fault dianosis
- Modeling sensor network data

40

# A general Bayes net

- Set of random variables

- Directed acyclic graph
  - Encodes independence assumptions

- CPTs

- Joint distribution:

$$P(X_1, \ldots, X_n) = \prod_{i=1}^{n} P\left(X_i \mid \mathbf{Pa}_{X_i}\right)$$

41

# How many parameters in a BN?

- Discrete variables $X_1, \ldots, X_n$
- Graph
  - Defines parents of $X_i$, $\mathbf{Pa}_{X_i}$
- CPTs – $P(X_i \mid \mathbf{Pa}_{Xi})$

42

# Another example

- Variables:
  - B – Burglar
  - E – Earthquake
  - A – Burglar alarm
  - N – Neighbor calls
  - R – Radio report
- Both burglars and earthquakes can set off the alarm
- If the alarm sounds, a neighbor may call
- An earthquake may be announced on the radio

**43**

# Another example – Building the BN

- B – Burglar
- E – Earthquake
- A – Burglar alarm
- N – Neighbor calls
- R – Radio report

**44**

# Independencies encoded in BN

- We said: All you need is the local Markov assumption
  - $(X_i \perp \text{NonDescendants}_{X_i} \mid \mathbf{Pa}_{X_i})$
- 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!!!

# Understanding independencies in BNs – BNs with 3 nodes

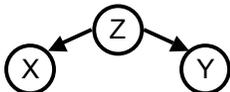**Local Markov Assumption:** A variable X is independent of its non-descendants given its parents
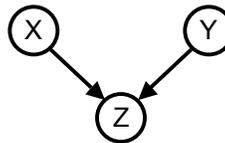
Indirect causal effect:



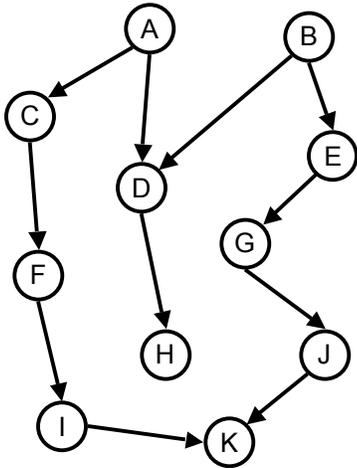Indirect evidential effect:



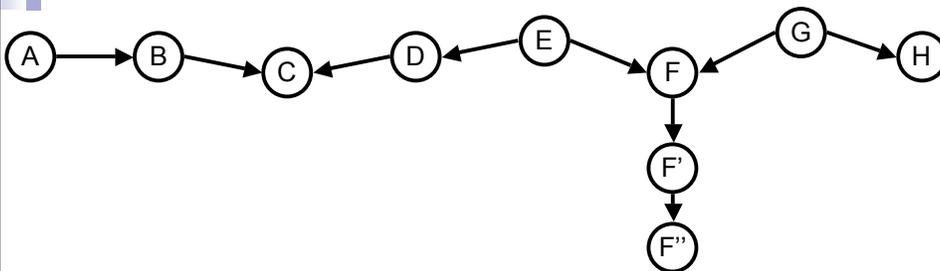Common cause:



Common effect:

# Understanding independencies in BNs – Some examples



©2005-2007 Carlos Guestrin

47

# An active trail – Example



**When are A and H independent?**

©2005-2007 Carlos Guestrin
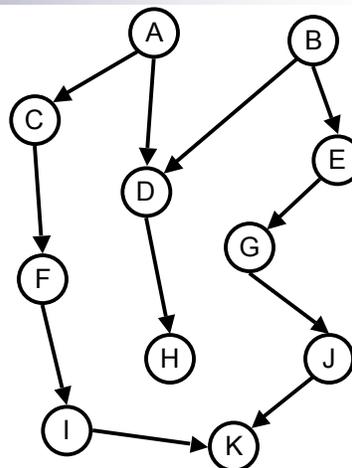
48

24

# Active trails formalized

- A path $X_1 - X_2 - \cdots - X_k$ is an **active trail** when variables $\boldsymbol{O} \subseteq \{X_1, \ldots, X_n\}$ 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** ($X_i \notin \boldsymbol{O}$)

  - $X_{i-1} \leftarrow X_i \leftarrow X_{i+1}$, and $X_i$ is **not observed** ($X_i \notin \boldsymbol{O}$)

  - $X_{i-1} \leftarrow X_i \rightarrow X_{i+1}$, and $X_i$ is **not observed** ($X_i \notin \boldsymbol{O}$)

  - $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$, and $X_i$ **is observed** ($X_i \in \boldsymbol{O}$), or **one of its descendents**

49

# Active trails and independence?

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

50

25

# The BN Representation Theorem

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

**Obtain** →

**Joint probability distribution:**

$$P(X_1, \ldots, X_n) = \prod_{i=1}^{n} P\left(X_i \mid \mathbf{Pa}_{X_i}\right)$$

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

**If joint probability distribution:**
$$P(X_1, \ldots, X_n) = \prod_{i=1}^{n} P\left(X_i \mid \mathbf{Pa}_{X_i}\right)$$

**Obtain** →

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

**Important because:**
**Read independencies of *P* from BN structure *G***

51

---

# "Simpler" BNs

- A distribution can be represented by many BNs:

- Simpler BN, requires fewer parameters

52

# Learning Bayes nets

|  | Known structure | Unknown structure |
|---|---|---|
| Fully observable data |  |  |
| Missing data |  |  |



Data

$\mathbf{x}^{(1)}$

...

$\mathbf{x}^{(m)}$

**structure** $+$ CPTs – $P(X_i | \mathbf{Pa}_{Xi})$ **parameters**

53

---

# Learning the CPTs

Data

$\mathbf{x}^{(1)}$

...

$\mathbf{x}^{(m)}$

For each discrete variable $X_i$



MLE: $P(X_i = x_i \mid X_j = x_j) = \dfrac{\mathsf{Count}(X_i = x_i, X_j = x_j)}{\mathsf{Count}(X_j = x_j)}$

54

# Queries in Bayes nets

- Given BN, find:
  - Probability of X given some evidence, $P(X|e)$

  - Most probable explanation, $\max_{x_1,\ldots,x_n} P(x_1,\ldots,x_n \mid e)$

  - Most informative query

- Learn more about these next class

55

# What you need to know

- Bayesian networks
  - A compact **representation** for large probability distributions
  - Not an algorithm
- Semantics of a BN
  - Conditional independence assumptions
- Representation
  - Variables
  - Graph
  - CPTs
- Why BNs are useful
- Learning CPTs from fully observable data
- Play with applet!!! ☺

56

# Acknowledgements

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

57