

# Bayes optimal classifier

## Naïve Bayes

Machine Learning – CSE446

Carlos Guestrin

University of Washington

November 18, 2014

©Carlos Guestrin 2005-2014

1

## Classification

- **Learn:**  $h: \mathbf{X} \mapsto Y$

- $\mathbf{X}$  – features
- $Y$  – target classes

- Suppose you know  $P(Y|\mathbf{X})$  exactly, how should you classify?

- Bayes optimal classifier:

©Carlos Guestrin 2005-2014

2

# Bayes Rule

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

Which is shorthand for:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{P(X = x_j)}$$

©Carlos Guestrin 2005-2014

3

# How hard is it to learn the optimal classifier?

■ Data =

Sky	Temp	Humid	Wind	Water	Forecst	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

■ How do we represent these? How many parameters?

□ Prior,  $P(Y)$ :

■ Suppose  $Y$  is composed of  $k$  classes

□ Likelihood,  $P(\mathbf{X}|Y)$ :

■ Suppose  $\mathbf{X}$  is composed of  $d$  binary features

■ **Complex model ! High variance with limited data!!!**

©Carlos Guestrin 2005-2014

4

## Conditional Independence

- X is **conditionally independent** of Y given Z, if the probability distribution governing X is independent of the value of Y, given the value of Z  
 $(\forall i, j, k) P(X = i | Y = j, Z = k) = P(X = i | Z = k)$

- e.g.,  $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$

- Equivalent to:

$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$

©Carlos Guestrin 2005-2014

5

## What if features are independent?

- Predict Thunder
- From two **conditionally Independent** features
  - Lightening
  - Rain

©Carlos Guestrin 2005-2014

6

## The Naïve Bayes assumption

- Naïve Bayes assumption:

- Features are independent given class:

$$\begin{aligned}P(X_1, X_2|Y) &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y)\end{aligned}$$

- More generally:

$$P(X_1 \dots X_d | Y) = \prod_i P(X_i | Y)$$

- How many parameters now?

- Suppose  $\mathbf{X}$  is composed of  $d$  binary features

©Carlos Guestrin 2005-2014

7

## The Naïve Bayes Classifier

- Given:

- Prior  $P(Y)$
- $d$  conditionally independent features  $\mathbf{X}$  given the class  $Y$
- For each  $X_i$ , we have likelihood  $P(X_i|Y)$

- Decision rule:

$$\begin{aligned}y^* = h_{NB}(\mathbf{x}) &= \arg \max_y P(y)P(x_1, \dots, x_d | y) \\ &= \arg \max_y P(y) \prod_i P(x_i | y)\end{aligned}$$

- If assumption holds, NB is optimal classifier!

©Carlos Guestrin 2005-2014

8

## MLE for the parameters of NB

- Given dataset
  - $\text{Count}(A=a, B=b) ==$  number of examples where  $A=a$  and  $B=b$
- MLE for NB, simply:
  - Prior:  $P(Y=y) =$
  
  - Likelihood:  $P(X_i=x_i|Y=y) =$

©Carlos Guestrin 2005-2014

9

## Subtleties of NB classifier 1 – Violating the NB assumption

- Usually, features are not conditionally independent:

$$P(X_1 \dots X_d | Y) \neq \prod_i P(X_i | Y)$$

- Actual probabilities  $P(Y|\mathbf{X})$  often biased towards 0 or 1
- Nonetheless, NB is the single most used classifier out there
  - NB often performs well, even when assumption is violated
  - [Domingos & Pazzani '96] discuss some conditions for good performance

©Carlos Guestrin 2005-2014

10

## Subtleties of NB classifier 2 – Insufficient training data

- What if you never see a training instance where  $X_1=a$  when  $Y=b$ ?
  - e.g.,  $Y=\{\text{SpamEmail}\}$ ,  $X_1=\{\text{'Enlargement'}\}$
  - $P(X_1=a \mid Y=b) = 0$
- Thus, no matter what the values  $X_2, \dots, X_d$  take:
  - $P(Y=b \mid X_1=a, X_2, \dots, X_d) = 0$
- “Solution”: smoothing
  - Add “fake” counts, usually uniformly distributed
  - Equivalent to Bayesian Learning

©Carlos Guestrin 2005-2014

11

## Text classification

- Classify e-mails
  - $Y = \{\text{Spam, NotSpam}\}$
- Classify news articles
  - $Y = \{\text{what is the topic of the article?}\}$
- Classify webpages
  - $Y = \{\text{Student, professor, project, ...}\}$
- What about the features  $\mathbf{X}$ ?
  - The text!

©Carlos Guestrin 2005-2014

12

## Features $\mathbf{X}$ are entire document – $X_i$ for $i^{\text{th}}$ word in article

### Article from rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e  
From: xxx@yyy.zzz.edu (John Doe)  
Subject: Re: This year's biggest and worst (opinic  
Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudehy is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some things in Toronto decided

13

## NB for Text classification

- $P(\mathbf{X}|Y)$  is huge!!!
  - Article at least 1000 words,  $\mathbf{X}=\{X_1, \dots, X_{1000}\}$
  - $X_i$  represents  $i^{\text{th}}$  word in document, i.e., the domain of  $X_i$  is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.
- NB assumption helps a lot!!!
  - $P(X_i=x_i|Y=y)$  is just the probability of observing word  $x_i$  in a document on topic  $y$

$$h_{NB}(\mathbf{x}) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

# Bag of words model

- Typical additional assumption – **Position in document doesn't matter**:  $P(X_i=x_i|Y=y) = P(X_k=x_i|Y=y)$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

When the lecture is over, remember to wake up the person sitting next to you in the lecture room.

# Bag of words model

- Typical additional assumption – **Position in document doesn't matter**:  $P(X_i=x_i|Y=y) = P(X_k=x_i|Y=y)$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

in is lecture lecture next over person remember room  
sitting the the the to to up wake when you



# Bag of Words Approach

The image shows a screenshot of a TOTAL website page titled "all about the company". The page contains several paragraphs of text. To the right of the screenshot is a table representing the word counts for various words in the document.

aardvark	0
about	2
all	2
Africa	1
apple	0
anxious	0
...	
gas	1
...	
oil	1
...	
Zaire	0

©Carlos Guestrin 2005-2014 17

## NB with Bag of Words for text classification

- Learning phase:
  - Prior  $P(Y)$ 
    - Count how many documents you have from each topic (+ prior)
  - $P(X_i|Y)$ 
    - For each topic, count how many times you saw word in documents of this topic (+ prior)
- Test phase:
  - For each document
    - Use naïve Bayes decision rule

$$h_{NB}(x) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

# Twenty News Groups results

Given 1000 training documents from each group  
Learn to classify new documents according to  
which newsgroup it came from

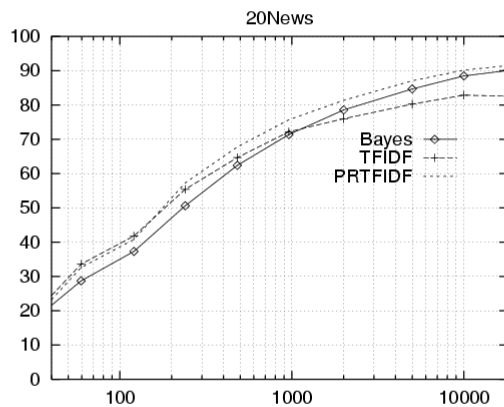
comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	
talk.politics.guns	

Naive Bayes: 89% classification accuracy

©Carlos Guestrin 2005-2014

19

# Learning curve for Twenty News Groups



Accuracy vs. Training set size (1/3 withheld for test)

©Carlos Guestrin 2005-2014

20

# Bayesian Networks – Representation

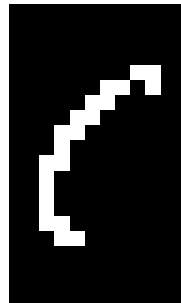
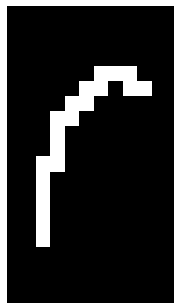
Machine Learning – CSE446  
Carlos Guestrin  
University of Washington

November 18, 2014

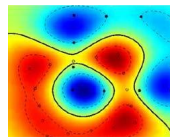
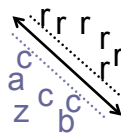
©Carlos Guestrin 2005-2014

21

## Handwriting recognition



Character recognition, e.g., kernel SVMs



©Carlos Guestrin 2005-2014

22

# Webpage classification



Company home page

VS

Personal home page

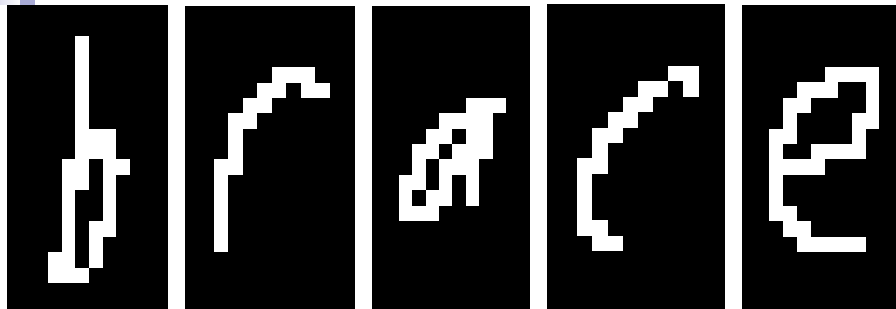
VS

University home page

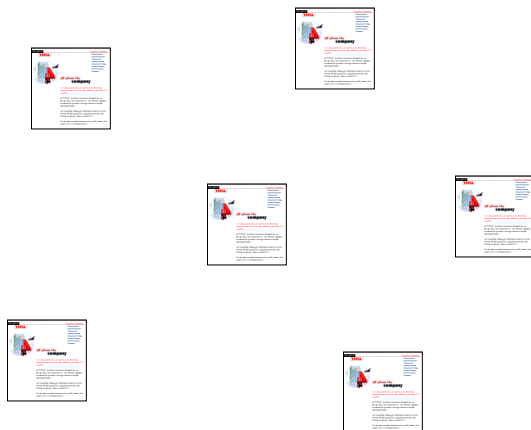
VS

...

# Handwriting recognition 2



## Webpage classification 2



©Carlos Guestrin 2005-2014

25

## Today – Bayesian networks

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

©Carlos Guestrin 2005-2014

26

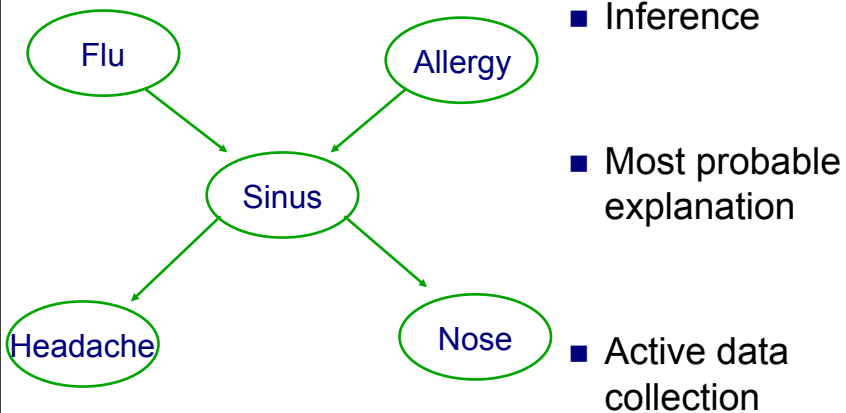
# Causal structure

- 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?

©Carlos Guestrin 2005-2014

27

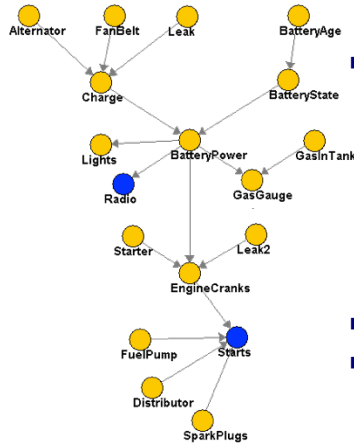
# Possible queries



©Carlos Guestrin 2005-2014

28

# Car starts BN

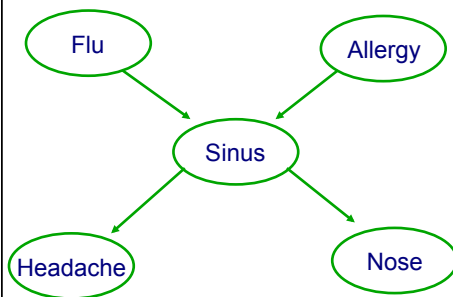


- 18 binary attributes
- Inference
  - $P(\text{BatteryAge} | \text{Starts}=f)$
- $2^{16}$  terms, why so fast?
- Not impressed?
  - HailFinder BN – more than  $3^{54} = 58149737003040059690390169$  terms

©Carlos Guestrin 2005-2014

29

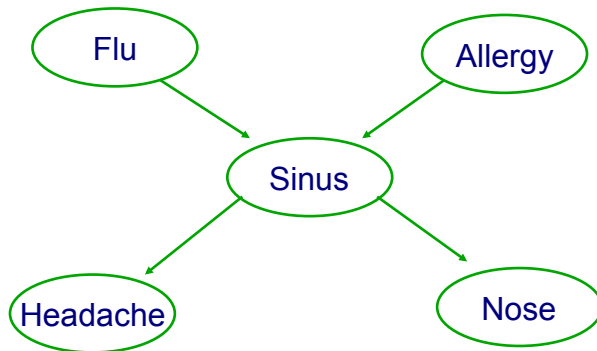
# Factored joint distribution - Preview



©Carlos Guestrin 2005-2014

30

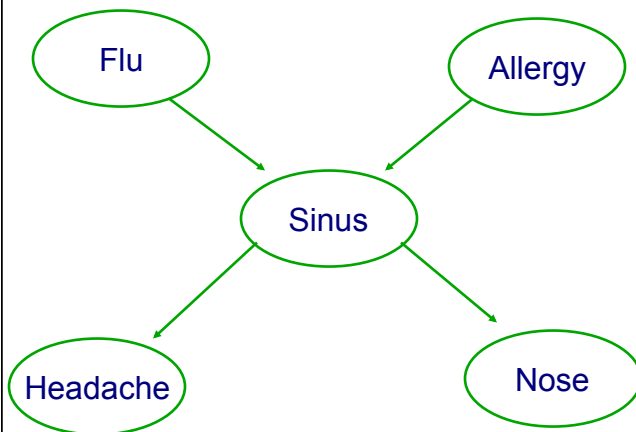
# What about probabilities? Conditional probability tables (CPTs)



©Carlos Guestrin 2005-2014

31

# Number of parameters

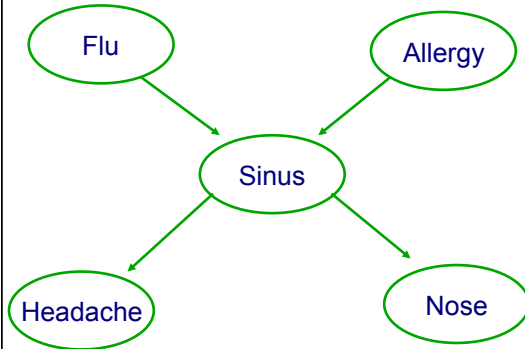


©Carlos Guestrin 2005-2014

32



## Key: Independence assumptions



Knowing sinus separates the variables from each other

©Carlos Guestrin 2005-2014

33

## (Marginal) Independence

- Flu and Allergy are (marginally) independent

Flu = t	
Flu = f	

Allergy = t	
Allergy = f	

	Flu = t	Flu = f
Allergy = t		
Allergy = f		

©Carlos Guestrin 2005-2014

34

## Marginally independent random variables

- **Sets** of variables  $\mathbf{X}, \mathbf{Y}$
- $\mathbf{X}$  is independent of  $\mathbf{Y}$  if
  - $P(\mathbf{X}=\mathbf{x} \perp \mathbf{Y}=\mathbf{y}), \forall \mathbf{x} \in \text{Val}(\mathbf{X}), \mathbf{y} \in \text{Val}(\mathbf{Y})$
- Shorthand:
  - **Marginal independence:**  $P(\mathbf{X} \perp \mathbf{Y})$
- **Proposition:**  $P$  satisfies  $(\mathbf{X} \perp \mathbf{Y})$  if and only if
  - $P(\mathbf{X}, \mathbf{Y}) = P(\mathbf{X}) P(\mathbf{Y})$

©Carlos Guestrin 2005-2014

35

## Conditional independence

- Flu and Headache are not (marginally) independent
- Flu and Headache are independent given Sinus infection
- More Generally:

©Carlos Guestrin 2005-2014

36

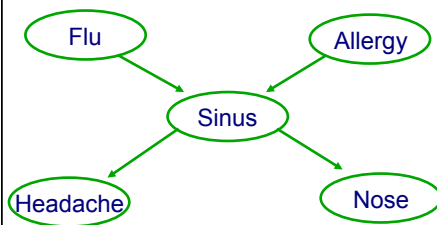
## Conditionally independent random variables

- **Sets** of variables **X**, **Y**, **Z**
- X is independent of Y given Z if
  - $P F (X=x \perp Y=y | Z=z), \forall x \in \text{Val}(X), y \in \text{Val}(Y), z \in \text{Val}(Z)$
- Shorthand:
  - **Conditional independence:**  $P F (X \perp Y | Z)$
  - For  $P F (X \perp Y | \emptyset)$ , write  $P F (X \perp Y)$
- **Proposition:**  $P$  satisfies  $(X \perp Y | Z)$  if and only if
  - $P(X, Y | Z) = P(X | Z) P(Y | Z)$

©Carlos Guestrin 2005-2014

37

## The independence assumption



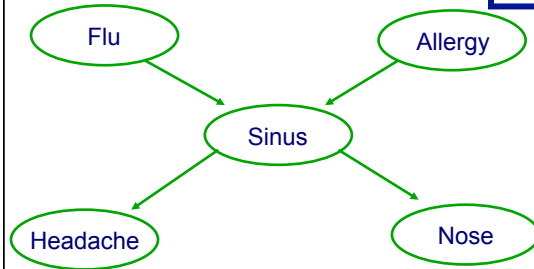
**Local Markov Assumption:**  
A variable X is independent of its non-descendants given its parents

©Carlos Guestrin 2005-2014

38

## Explaining away

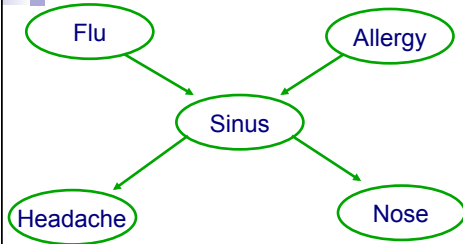
**Local Markov Assumption:**  
A variable X is independent of its non-descendants given its parents



## Naïve Bayes revisited

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

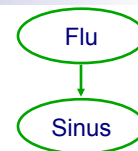
## Joint distribution



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

## The chain rule of probabilities

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



- More generally:

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

# Chain rule & Joint distribution

```

    graph TD
      Flu((Flu)) --> Sinus((Sinus))
      Allergy((Allergy)) --> Sinus((Sinus))
      Sinus((Sinus)) --> Headache((Headache))
      Sinus((Sinus)) --> Nose((Nose))
  
```

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

©Carlos Guestrin 2005-2014 43

# 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:

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

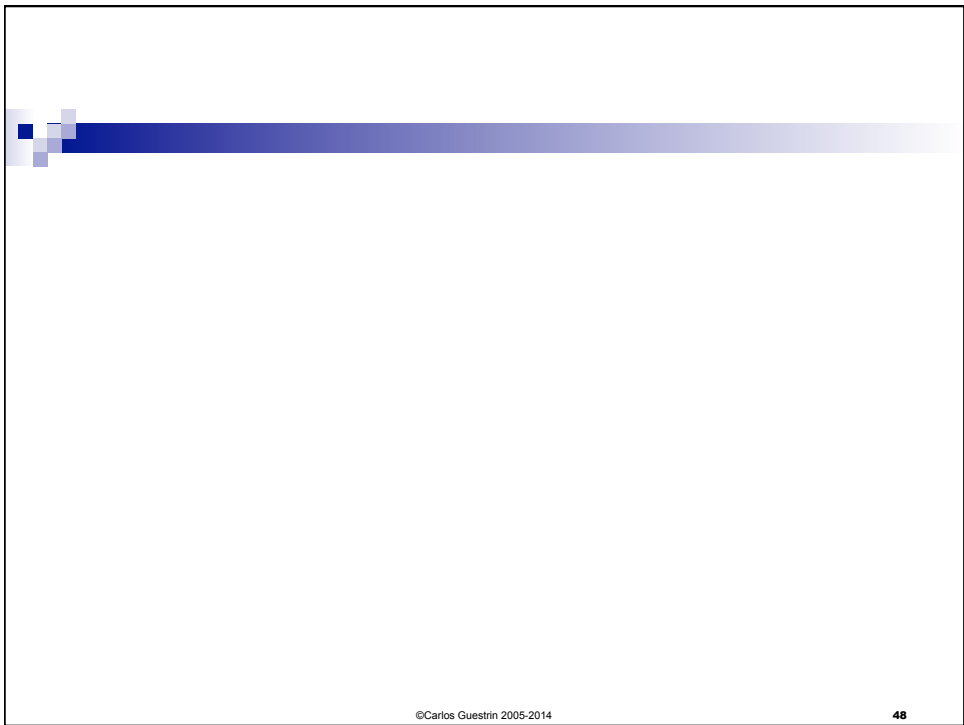
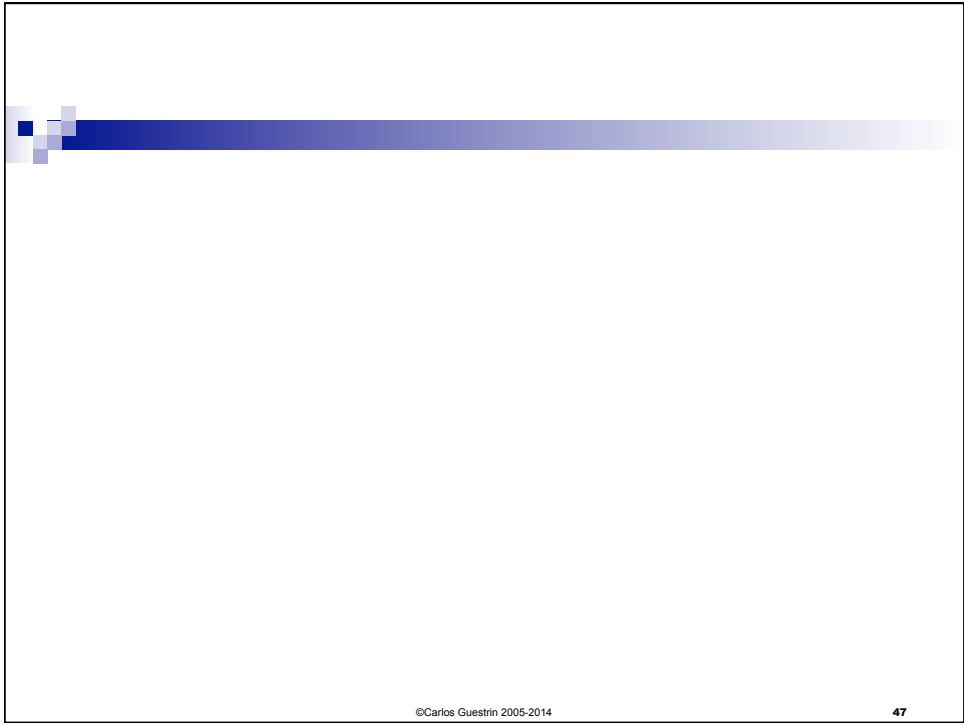
©Carlos Guestrin 2005-2014 44

# Two (trivial) special cases

Edgeless graph

Fully-connected  
graph







## Feature selection

- Want to learn  $f: \mathbf{X} \rightarrow Y$ 
  - $\mathbf{X} = \langle X_1, \dots, X_n \rangle$
  - but some features are more important than others
- **Approach:** select subset of features to be used by learning algorithm
  - **Score** each feature (or sets of features)
  - **Select** set of features with best score

©Carlos Guestrin 2005-2014

## Simple greedy **forward** feature selection algorithm

- Pick a dictionary of features
  - e.g., polynomials for linear regression
- Greedy heuristic:
  - Start from empty (or simple) set of features  $F_0 = \emptyset$
  - Run learning algorithm for current set of features  $F_t$ 
    - Obtain  $h_t$
  - Select **next best feature**  $X_i$ 
    - e.g.,  $X_j$  that results in lowest cross-validation error learner when learning with  $F_t \cup \{X_j\}$
  - $F_{t+1} \leftarrow F_t \cup \{X_i\}$
  - Recurse

©Carlos Guestrin 2005-2014

## Simple greedy **backward** feature selection algorithm

- Pick a dictionary of features
  - e.g., polynomials for linear regression
- Greedy heuristic:
  - Start from all features  $F_0 = F$
  - Run learning algorithm for current set of features  $F_t$ 
    - Obtain  $h_t$
  - Select **next worst feature**  $X_i$ 
    - e.g.,  $X_j$  that results in lowest cross-validation error learner when learning with  $F_t - \{X_j\}$
  - $F_{t+1} \leftarrow F_t - \{X_i\}$
  - Recurse

©Carlos Guestrin 2005-2014

## Impact of feature selection on classification of fMRI data [Pereira et al. '05]

Accuracy classifying category of word read by subject

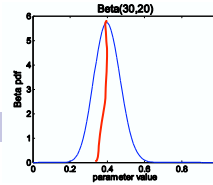
#voxels	mean	subjects								
		233B	329B	332B	424B	474B	496B	77B	86B	
50	0.735	0.783	0.817	0.55	0.783	0.75	0.8	0.65	0.75	
100	0.742	0.767	0.8	0.533	0.817	0.85	0.783	0.6	0.783	
200	0.737	0.783	0.783	0.517	0.817	0.883	0.75	0.583	0.783	
<b>300</b>	<b>0.75</b>	<b>0.8</b>	<b>0.817</b>	<b>0.567</b>	<b>0.833</b>	<b>0.883</b>	<b>0.75</b>	<b>0.583</b>	<b>0.767</b>	
400	0.742	0.8	0.783	0.583	0.85	0.833	0.75	0.583	0.75	
800	0.735	0.833	0.817	0.567	0.833	0.833	0.7	0.55	0.75	
1600	0.698	0.8	0.817	0.45	0.783	0.833	0.633	0.5	0.75	
all (~2500)	0.638	0.767	0.767	0.25	0.75	0.833	0.567	0.433	0.733	

Table 1: Average accuracy across all pairs of categories, restricting the procedure to use a certain number of voxels for each subject. The highlighted line corresponds to the best mean accuracy, obtained using 300 voxels.

Voxels scored by p-value of regression to predict voxel value from the task

©Carlos Guestrin 2005-2014

## MAP for Beta distribution



$$P(\theta | \mathcal{D}) = \frac{\theta^{\beta_H + \alpha_H - 1} (1 - \theta)^{\beta_T + \alpha_T - 1}}{B(\beta_H + \alpha_H, \beta_T + \alpha_T)} \sim \text{Beta}(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

- MAP: use most likely parameter:

$$\hat{\theta} = \arg \max_{\theta} P(\theta | \mathcal{D}) = \frac{\beta_H + \alpha_H - 1}{\beta_H + \alpha_H + \beta_T + \alpha_T - 2}$$

$\alpha_H = 3$   
 $\alpha_T = 2$   
 $\beta_H, \beta_T$  extra data

- Beta prior equivalent to extra thumbtack flips
- As  $N \rightarrow 1$ , prior is “forgotten”
- **But, for small sample size, prior is important!**

©Carlos Guestrin 2005-2014

53

## Bayesian learning for NB parameters – a.k.a. smoothing

- Dataset of  $N$  examples
- Prior
  - “distribution”  $Q(X_i, Y), Q(Y)$
  - $m$  “virtual” examples
- MAP estimate
  - $P(X_i | Y)$
- **Now, even if you never observe a feature/class, posterior probability never zero**

©Carlos Guestrin 2005-2014

54

# Properties of independence

- **Symmetry:**

- $(X \perp Y | Z) \Rightarrow (Y \perp X | Z)$

- **Decomposition:**

- $(X \perp Y, W | Z) \Rightarrow (X \perp Y | Z)$

- **Weak union:**

- $(X \perp Y, W | Z) \Rightarrow (X \perp Y | Z, W)$

- **Contraction:**

- $(X \perp W | Y, Z) \& (X \perp Y | Z) \Rightarrow (X \perp Y, W | Z)$

- **Intersection:**

- $(X \perp Y | W, Z) \& (X \perp W | Y, Z) \Rightarrow (X \perp Y, W | Z)$

- Only for positive distributions!

- $P(\alpha) > 0, \forall \alpha, \alpha > 0;$