

Naïve Bayes (Continued) Naïve Bayes with Continuous (variables) Logistic Regression

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

January 29th, 2007

©2005-2007 Carlos Guestrin

1

Announcements

- Recitations stay on Thursdays
 - 5-6:30pm in Wean 5409
 - This week: Naïve Bayes & Logistic Regression
- **Extension** for the first homework:
 - **Due Wed. Feb 8th** beginning of class
 - Mitchell's chapter is most useful reading
- Go to the AI seminar:
 - Tuesdays 3:30pm, Wean 5409
 - <http://www.cs.cmu.edu/~aiseminar/>
 - This week's seminar very relevant to what we are covering in class

*ignore, unless
you are taking
the class
last year...*

©2005-2007 Carlos Guestrin

2

Optimal classification $P(Y|X)$

- **Theorem:** Bayes classifier h_{Bayes} is optimal! *if you know $P(Y|X)$ exactly*

$$h_{\text{Bayes}}(x) = y^*(x) \in \arg \max_y P(Y=y|X=x)$$

□ That is $\text{error}_{\text{true}}(h_{\text{Bayes}}) \leq \text{error}_{\text{true}}(h), \forall h(x)$

- **Proof:** $p(\text{error}) = \int_x p(\text{error}|x)p(x)dx$

$$\min_h p(\text{error}|x) = \begin{cases} P(Y=f|x), & h(x)=t \\ P(Y=t|x), & h(x)=f \end{cases} \begin{aligned} &= \max \text{ prob get right} \\ &= \text{guess answer with highest prob} \\ &= \arg \max_y P(Y=y|X=x) \end{aligned}$$

©2005-2007 Carlos Guestrin

3

How hard is it to learn the optimal classifier?

- Data =

Sky	Temp	Humid	Wind	Water	Forecast	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, $P(Y)$: $k-1$ ← *Probably OK*
 - Suppose Y is composed of k classes

- Likelihood, $P(X|Y)$: $k(2^n - 1)$ ← *get you into trouble*
 - Suppose X is composed of n binary features

$$\# \text{param}(P(X|Y=y)) = 2^n - 1$$

↑
for each y

- Complex model → High variance with limited data!!!

©2005-2007 Carlos Guestrin

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

$T \perp R | L$ = Thunder independent of rain given lightning

- Equivalent to:

$X \perp Y | Z$

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

©2005-2007 Carlos Guestrin

5

The Naïve Bayes assumption

- Naïve Bayes assumption:

$X_1 \perp X_2 | Y$

- Features are independent given class:

likelihood

$$P(X_1, X_2 | Y) = P(X_1 | X_2, Y)P(X_2 | Y) \\ = P(X_1 | Y)P(X_2 | Y)$$

- More generally:

$\forall i, X_i \perp \{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\} | Y$

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

- How many parameters now?

- Suppose X is composed of n binary features

full table \rightarrow $2 \cdot (2^n - 1)$ parameters

with NB assumption $n \cdot K \cdot (2 - 1) = nK$

$Z = m$ values

$P(Z) = \frac{1}{m} P(Z)$

$\sum_j P(Z=j) = 1$

Saves one parameter

©2005-2007 Carlos Guestrin

6

The Naïve Bayes Classifier

■ Given:

- Prior $P(Y)$ *0.8* *0.2*
- n conditionally independent features X given the class Y
- For each X_i , we have likelihood $P(X_i|Y)$ *works images*

■ Decision rule:

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

compute for each feature prior likelihood

- If assumption holds, NB is optimal classifier!

©2005-2007 Carlos Guestrin

7

MLE for the parameters of NB

■ Given dataset D

- Count(A=a, B=b) ← number of examples where A=a and B=b

■ MLE for NB, simply:

- Prior: $P(Y=y) = \frac{\text{Count}(Y=y)}{|D|}$

- Likelihood: $P(X_i=x_i|Y_i=y_i) = \frac{\text{Count}(X_i=x_i, Y_i=y_i)}{\text{Count}(Y_i=y_i)}$

*proof path:
write likelihood of data
take log
take derivative
set to zero.
;*

©2005-2007 Carlos Guestrin

8

Subtleties of NB classifier 1 – Violating the NB assumption

- Usually, features are not conditionally independent:

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

- Actual probabilities $P(Y|X)$ 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

©2005-2007 Carlos Guestrin

9

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., $Y=\{\text{SpamEmail}\}$, $X_1=\{\text{'Enlargement'}\}$
- $P(X_1=a | Y=b) = 0$

- Thus, no matter what the values X_2, \dots, X_n take:

- $P(Y=b | X_1=a, X_2, \dots, X_n) = 0$

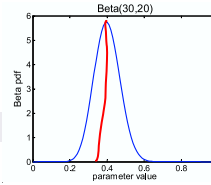
$$P(Y=b) \cdot P(X_1=a | Y=b) \cdot \prod_{i=2}^n P(X_i | Y=b) = 0$$

- What now???

©2005-2007 Carlos Guestrin

10

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:

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

from prior

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

©2005-2007 Carlos Guestrin

11

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

- Dataset of N examples *(assume Dirichlet prior)*
- Prior *(generalization of Beta for multinomials)*
 - "distribution" $Q(X_i, Y), Q(Y)$, typically $Q \sim \text{uniform}$
 - m "virtual" examples, $m = 1, 2, \dots, 5, \dots, 10$

- MAP estimate

$$P(X_i = a | Y) = \frac{\text{Count}(X_i = a, Y = b) + m \cdot Q(X_i = a, Y = b)}{\text{Count}(Y = b) + m \cdot Q(Y = b)}$$

e.g., if Q is uniform, $|Y| = k$, $|X_i| = \ell$

$$Q(Y = b) = \frac{1}{k} \quad Q(X_i = a, Y = b) = \frac{1}{k \cdot \ell}$$

- **Now, even if you never observe a feature/class, posterior probability never zero**

©2005-2007 Carlos Guestrin

12

Text classification

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

©2005-2007 Carlos Guestrin

13

Features X are entire document – X_i for i^{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 (opinion)
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 Hrudey 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 thugs in Toronto decided

$X_i = \{\text{Kings}\}$
 $|X_i| = 10,000 \rightarrow 109,000$

©2005-2007 Carlos Guestrin

14

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^{th} word in document, i.e., the domain of X_i is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.

explicitly $\#_{\text{params}}(P(\mathbf{X}|\mathbf{Y})) = K \cdot (10,000^{1000} - 1)$

■ 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**: $P(X_i=x_i|Y=y) = P(X_k=x_k|Y=y)$

- “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

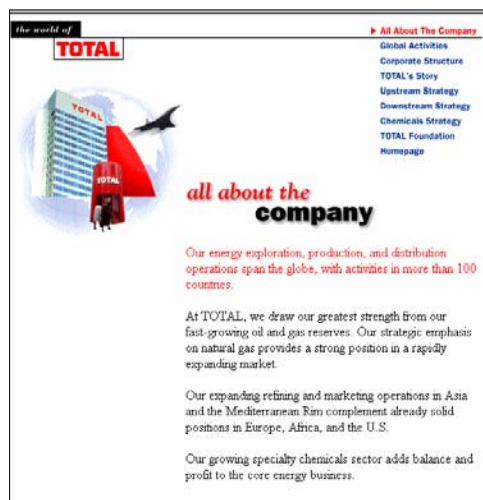
- Typical additional assumption – **Position in document doesn't matter**: $P(X_i=x_i|Y=y) = P(X_k=x_k|Y=y)$

- “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 next over person remember room
sitting the the the to to up wake when you

Bag of Words Approach



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

NB with Bag of Words for text classification

■ Learning phase:

□ Prior $P(Y)$

- Count how many documents you have from each topic (+ *smoothing* 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

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

©2005-2007 Carlos Guestrin

19

Twenty News Groups results

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

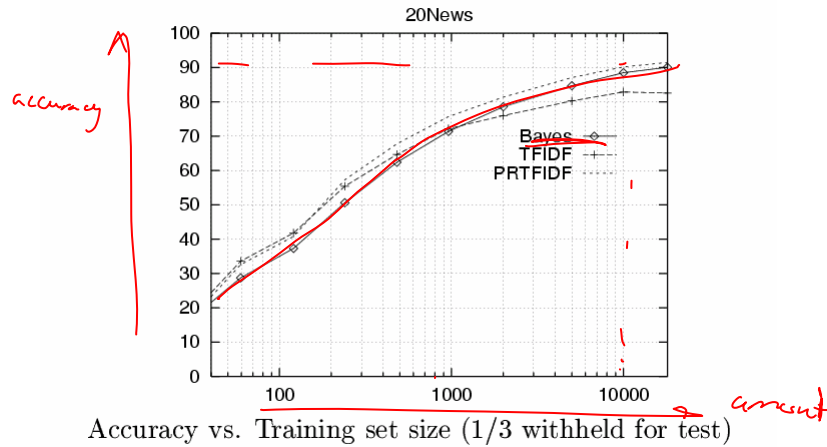
comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc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

Naive Bayes 89% classification accuracy

©2005-2007 Carlos Guestrin

20

Learning curve for Twenty News Groups



©2005-2007 Carlos Guestrin

21

What if we have continuous X_i ?

Eg., character recognition (X_i is i th pixel) $Y \in \{a, b, \dots, z\}$

$$P(X|Y) = P \left(\begin{matrix} \text{image} \\ Y \end{matrix} \right)$$

Handwritten red annotations: $x_i \in \{0, \dots, 255\}$ with an arrow pointing to a pixel in the image, and x_i with an arrow pointing to the i th pixel.



Gaussian Naïve Bayes (GNB):

$$P(X_i = x | Y = y_k) = \frac{1}{\sigma_{ik} \sqrt{2\pi}} e^{-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}}$$

Handwritten red annotations: $x_i \in (-\infty, +\infty)$ with an arrow pointing to the x variable; μ_{ik} = for feature i in class k with an arrow pointing to μ_{ik} ; σ_{ik} = for feature i in class k with an arrow pointing to σ_{ik} .

Sometimes assume variance

- is independent of Y (i.e., σ_i) \leftarrow does not depend class
- or independent of X_i (i.e., σ_k) \leftarrow does not depend on feature
- or both (i.e., σ) \leftarrow nothing

©2005-2007 Carlos Guestrin

22

Estimating Parameters: Y discrete, X_i continuous

Maximum likelihood estimates:

$$\hat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_{j=1}^{|D|} X_i^j \delta(Y^j = y_k)$$

*5th feature
kth class*

jth training
example

$\delta(x)=1$ if x true,
else 0

$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k) - 1} \sum_{j=1}^{|D|} (X_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k)$$

©2005-2007 Carlos Guestrin

23

Example: GNB for classifying mental states

[Mitchell et al.]

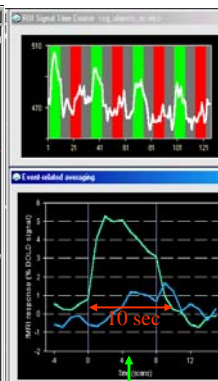
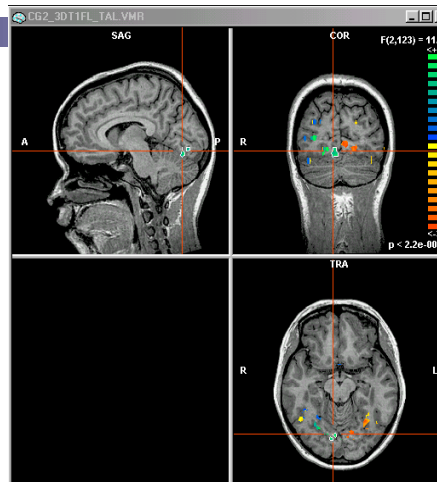
~1 mm resolution

~2 images per sec.

15,000 voxels/image

non-invasive, safe

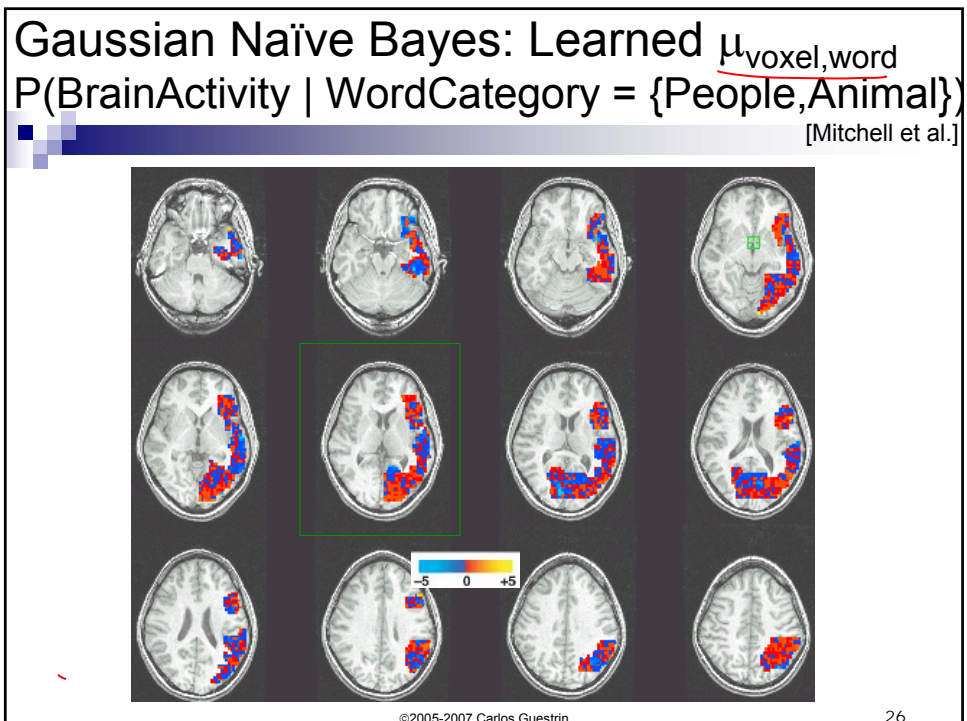
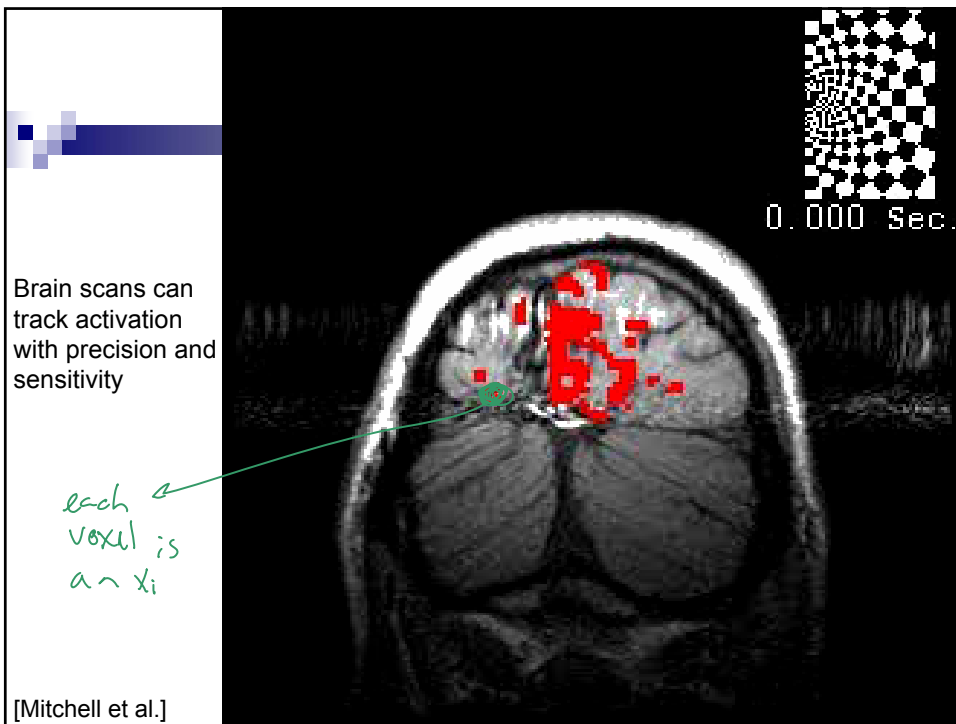
measures Blood
Oxygen Level
Dependent (BOLD)
response



Typical
impulse
response

©2005-2007 Carlos Guestrin

24



Learned Bayes Models – Means for $P(\text{BrainActivity} \mid \text{WordCategory})$

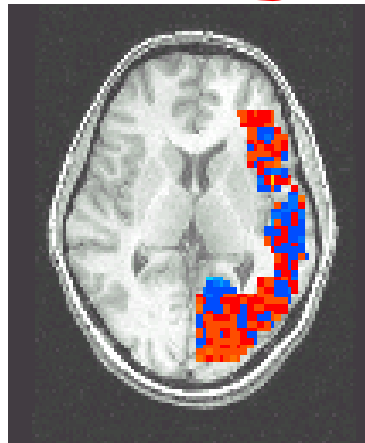
[Mitchell et al.]

Pairwise classification accuracy: 85%

People words



Animal words



©2005-2007 Carlos Guestrin

27

What you need to know about Naïve Bayes

- Types of learning problems
 - Learning is (just) function approximation!
- Optimal decision using Bayes Classifier $P(y|x)$
- Naïve Bayes classifier $P(x)$, $P(x|y)$, use Bayes rule... NB
 - What's the assumption
 - Why we use it
 - How do we learn it
 - Why is Bayesian estimation important
- Text classification
 - Bag of words model
- Gaussian NB
 - Features are still conditionally independent
 - Each feature has a Gaussian distribution given class

©2005-2007 Carlos Guestrin

28

Generative v. Discriminative classifiers – Intuition

Want to Learn: $h: X \mapsto Y$

- X – features
- Y – target classes

Bayes optimal classifier – $P(Y|X)$

Generative classifier, e.g., Naïve Bayes:

- Assume some functional form for $P(X|Y)$, $P(Y)$
- Estimate parameters of $P(X|Y)$, $P(Y)$ directly from training data
- Use Bayes rule to calculate $P(Y|X=x) \propto P(Y) \cdot P(X=x|Y)$
- This is a 'generative' model
- Indirect computation of $P(Y|X)$ through Bayes rule
- But, can generate a sample of the data, $P(X) = \sum_y P(y) P(X|y)$

Discriminative classifiers, e.g., Logistic Regression:

- Assume some functional form for $P(Y|X)$
- Estimate parameters of $P(Y|X)$ directly from training data
- This is the 'discriminative' model
- Directly learn $P(Y|X)$
- But cannot obtain a sample of the data, because $P(X)$ is not available

$P(Y)$
 $P(X|Y)$
multiple images
first sample
class
then sample
pixel values

if you have
 $h_{\text{Bayes}}(x) = \arg \max_y P(Y|X=x)$

NB ← table

if you want to classify have use

directly model $P(Y|X)$

discriminate classes, e.g.,
person v. animal

©2005-2007 Carlos Guestrin

29

Logistic Regression

Logistic function (or Sigmoid):

$$\frac{1}{1 + \exp(-z)}$$

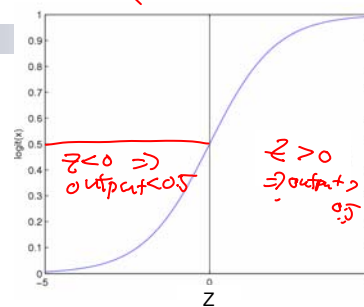
Learn $P(Y|X)$ directly!

- Assume a particular functional form
- Sigmoid applied to a linear function of the data:

$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

linear function of X , with parameters w

$$P(Y=0|X) = 1 - P(Y=1|X)$$



Features can be discrete or continuous!

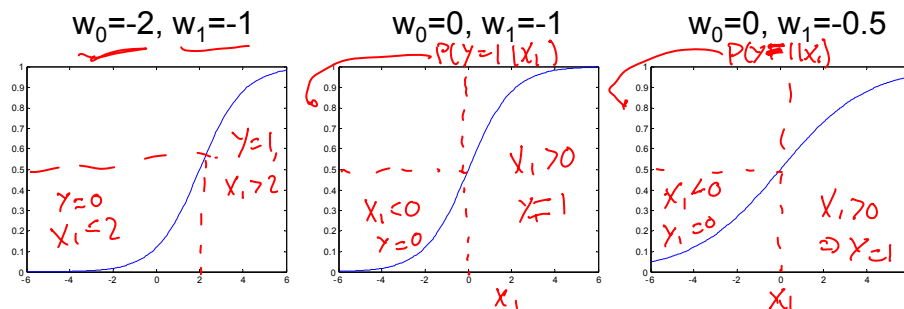
©2005-2007 Carlos Guestrin

30

Understanding the sigmoid

$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

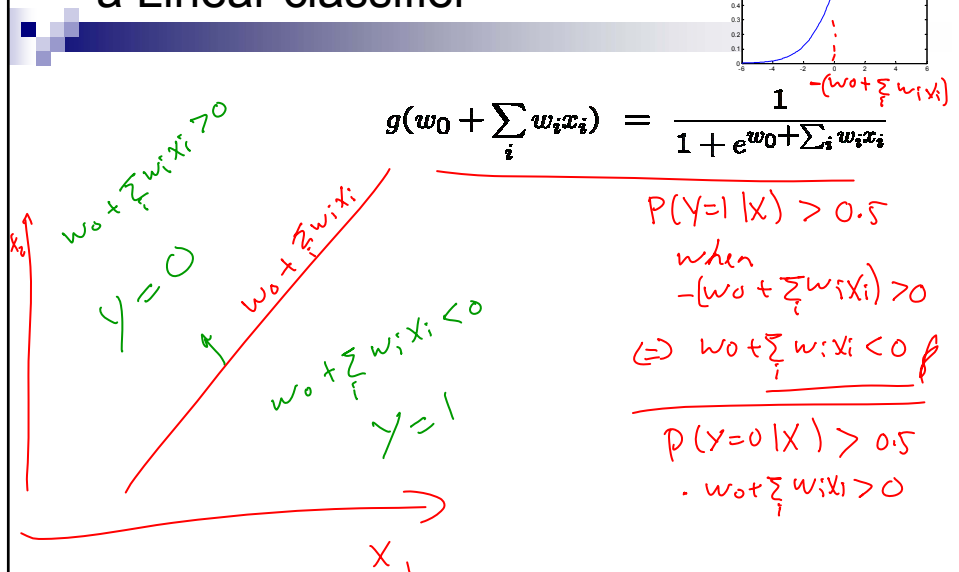
one feature x_1 :



©2005-2007 Carlos Guestrin

31

Logistic Regression – a Linear classifier



©2005-2007 Carlos Guestrin

32

Very convenient!

$$P(Y = 1|X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$P(Y = 0|X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$\frac{P(Y = 0|X)}{P(Y = 1|X)} = \exp(w_0 + \sum_i w_i X_i)$$

implies

$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$

linear
classification
rule!

©2005-2007 Carlos Guestrin

33

Logistic regression v. Naïve Bayes

- Consider learning $f: X \rightarrow Y$, where
 - X is a vector of real-valued features, $\langle X_1 \dots X_n \rangle$
 - Y is boolean
- Could use a Gaussian Naïve Bayes classifier
 - assume all X_i are conditionally independent given Y
 - model $P(X_i | Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i)$
 - model $P(Y)$ as Bernoulli($\theta, 1-\theta$)
- What does that imply about the form of $P(Y|X)$?

©2005-2007 Carlos Guestrin

34

Logistic regression v. Naïve Bayes

- Consider learning $f: X \rightarrow Y$, where
 - X is a vector of real-valued features, $\langle X_1 \dots X_n \rangle$
 - Y is boolean
- Could use a Gaussian Naïve Bayes classifier
 - assume all X_i are conditionally independent given Y
 - model $P(X_i | Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i)$
 - model $P(Y)$ as Bernoulli($\theta, 1-\theta$)

- What does that imply about the form of $P(Y|X)$?

$$P(Y = 1 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Cool!!!!

©2005-2007 Carlos Guestrin

35

Derive form for $P(Y|X)$ for continuous X_i

$$\begin{aligned}
 P(Y = 1 | X) &= \frac{P(Y = 1)P(X|Y = 1)}{P(Y = 1)P(X|Y = 1) + P(Y = 0)P(X|Y = 0)} \\
 &= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} \\
 &= \frac{1}{1 + \exp(\ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)})} \\
 &= \frac{1}{1 + \exp(\ln \frac{1-\theta}{\theta} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)})}
 \end{aligned}$$

©2005-2007 Carlos Guestrin

36

Ratio of class-conditional probabilities

$$\ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}$$

$$P(X_i = x | Y = y_k) = \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x - \mu_{ik})^2}{2\sigma_i^2}}$$

Derive form for $P(Y|X)$ for continuous X_i

$$P(Y=1|X) = \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)}$$

$$= \frac{1}{1 + \exp\left(\ln \frac{1-\theta}{\theta}\right) + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}}$$

$$\sum_i \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right)$$

$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

Gaussian Naïve Bayes v. Logistic Regression

■ **Set of Gaussian Naïve Bayes parameters**
(feature variance independent of class label)

■ **Set of Logistic Regression parameters**

- Representation equivalence
 - But only in a special case!!! (GNB with class-independent variances)
- But what's the difference???
- **LR makes no assumptions about $P(X|Y)$ in learning!!!**
- **Loss function!!!**
 - Optimize different functions → Obtain different solutions

©2005-2007 Carlos Guestrin

39

Logistic regression for more than 2 classes

- Logistic regression in more general case, where $Y \in \{Y_1 \dots Y_R\}$: learn $R-1$ sets of weights

©2005-2007 Carlos Guestrin

40

Logistic regression more generally

- Logistic regression in more general case, where $Y \in \{Y_1 \dots Y_R\}$: learn $R-1$ sets of weights

for $k < R$

$$P(Y = y_k | X) = \frac{\exp(w_{k0} + \sum_{i=1}^n w_{ki} X_i)}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} X_i)}$$

for $k=R$ (normalization, so no weights for this class)

$$P(Y = y_R | X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} X_i)}$$

Features can be discrete or continuous!

©2005-2007 Carlos Guestrin

41

Loss functions: Likelihood v. Conditional Likelihood

- Generative (Naïve Bayes) Loss function:

Data likelihood

$$\begin{aligned} \ln P(\mathcal{D} | \mathbf{w}) &= \sum_{j=1}^N \ln P(\mathbf{x}^j, y^j | \mathbf{w}) \\ &= \sum_{j=1}^N \ln P(y^j | \mathbf{x}^j, \mathbf{w}) + \sum_{j=1}^N \ln P(\mathbf{x}^j | \mathbf{w}) \end{aligned}$$

- Discriminative models cannot compute $P(\mathbf{x} | \mathbf{w})$!
- But, discriminative (logistic regression) loss function:

Conditional Data Likelihood

$$\ln P(\mathcal{D}_Y | \mathcal{D}_X, \mathbf{w}) = \sum_{j=1}^N \ln P(y^j | \mathbf{x}^j, \mathbf{w})$$

- Doesn't waste effort learning $P(X)$ – focuses on $P(Y|X)$ all that matters for classification

©2005-2007 Carlos Guestrin

42

Expressing Conditional Log Likelihood

$$l(\mathbf{w}) \equiv \sum_j \ln P(y^j | \mathbf{x}^j, \mathbf{w})$$

$$P(Y = 0 | \mathbf{X}, \mathbf{w}) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1 | \mathbf{X}, \mathbf{w}) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$l(\mathbf{w}) = \sum_j y^j \ln P(y^j = 1 | \mathbf{x}^j, \mathbf{w}) + (1 - y^j) \ln P(y^j = 0 | \mathbf{x}^j, \mathbf{w})$$

©2005-2007 Carlos Guestrin

43

Maximizing Conditional Log Likelihood

$$l(\mathbf{w}) \equiv \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w})$$

$$P(Y = 0 | \mathbf{X}, \mathbf{w}) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1 | \mathbf{X}, \mathbf{w}) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$= \sum_j y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))$$

Good news: $l(\mathbf{w})$ is concave function of $\mathbf{w} \rightarrow$ no locally optimal solutions

Bad news: no closed-form solution to maximize $l(\mathbf{w})$

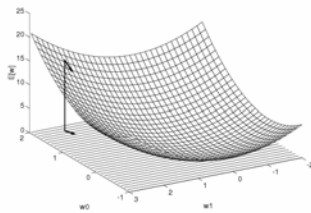
Good news: concave functions easy to optimize

©2005-2007 Carlos Guestrin

44

Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave
→ Find optimum with gradient ascent



Gradient: $\nabla_{\mathbf{w}} l(\mathbf{w}) = \left[\frac{\partial l(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial l(\mathbf{w})}{\partial w_n} \right]'$

Learning rate, $\eta > 0$

Update rule: $\Delta \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$

$$w_i \leftarrow w_i + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

- Gradient ascent is simplest of optimization approaches
 - e.g., Conjugate gradient ascent much better (see reading)

©2005-2007 Carlos Guestrin

45

Maximize Conditional Log Likelihood: Gradient ascent

$$l(\mathbf{w}) = \sum_j y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))$$

Gradient ascent algorithm: iterate until change $< \epsilon$

For all i , $w_i \leftarrow w_i + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w})]$
repeat

©2005-2007 Carlos Guestrin

46

That's all M(C)LE. How about MAP?

$$p(\mathbf{w} \mid Y, \mathbf{X}) \propto P(Y \mid \mathbf{X}, \mathbf{w})p(\mathbf{w})$$

- One common approach is to define priors on \mathbf{w}
 - Normal distribution, zero mean, identity covariance
 - “Pushes” parameters towards zero
- Corresponds to **Regularization**
 - Helps avoid very large weights and overfitting
 - Explore this in your homework
 - More on this later in the semester

- MAP estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[p(\mathbf{w}) \prod_{j=1}^N P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

©2005-2007 Carlos Guestrin

47

Gradient of M(C)AP

$$\frac{\partial}{\partial w_i} \ln \left[p(\mathbf{w}) \prod_{j=1}^N P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right] \qquad p(\mathbf{w}) = \prod_i \frac{1}{\kappa \sqrt{2\pi}} e^{\frac{-w_i^2}{2\kappa^2}}$$

©2005-2007 Carlos Guestrin

48

MLE vs MAP

- Maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[\prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right]$$

$$w_i \leftarrow w_i + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w})]$$

- Maximum conditional a posteriori estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[p(\mathbf{w}) \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right]$$

$$w_i \leftarrow w_i + \eta \left\{ -\lambda w_i + \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w})] \right\}$$

©2005-2007 Carlos Guestrin

49

What you should know about Logistic Regression (LR)

- Gaussian Naïve Bayes with class-independent variances representationally equivalent to LR
 - Solution differs because of objective (loss) function
- In general, NB and LR make different assumptions
 - NB: Features independent given class \rightarrow assumption on $P(\mathbf{X}|Y)$
 - LR: Functional form of $P(Y|\mathbf{X})$, no assumption on $P(\mathbf{X}|Y)$
- LR is a linear classifier
 - decision rule is a hyperplane
- LR optimized by conditional likelihood
 - no closed-form solution
 - concave \rightarrow global optimum with gradient ascent
 - Maximum conditional a posteriori corresponds to regularization

©2005-2007 Carlos Guestrin

50

Acknowledgements

- Some of the material in the presentation is courtesy of Tom Mitchell