# Naïve Bayes, Gaussian Distributions, Practical Applications

#### Required reading:

 Mitchell draft chapter, sections 1 and 2. (available on class website)

Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University



Ground Hog's Day, 2009

## Overview

#### Recently:

- learn P(Y|X) instead of deterministic f:X→Y
- Bayes rule
- MLE and MAP estimates for parameters of P
- Conditional independence
- classification with Naïve Bayes

#### Today:

- Text classification with Naïve bayes
- Gaussian distributions for continuous X
- Gaussian Naïve Bayes classifier
- Image classification with Naïve bayes

## Learning to classify text documents

- Classify which emails are spam?
- Classify which emails promise an attachment?
- Classify which web pages are student home pages?

How shall we represent text documents for Naïve Bayes?

#### Article from rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e

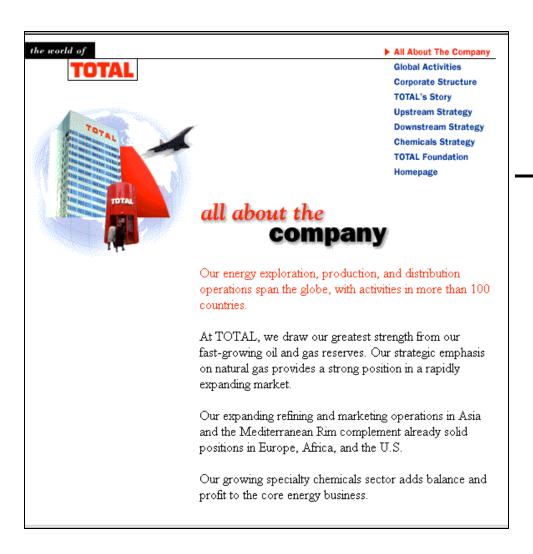
From: xxx@yyy.zzz.edu (John Doe)

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 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

## Baseline: Bag of Words Approach



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

## Naïve Bayes in a Nutshell

Bayes rule:  $P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) P(X_1 \dots X_n | Y = y_k)}{\sum_j P(Y = y_j) P(X_1 \dots X_n | Y = y_j)}$ Sport  $\in \{\text{hockey, bb-}\}$ 

Assuming conditional independence among X<sub>i</sub>'s:

$$P(Y = y_k | X_1 ... X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

So, classification rule for  $X^{new} = \langle X_1, ..., X_n \rangle$  is:

### Learning to Classify Text

Target concept  $Interesting?: Document \rightarrow \{+, -\}$ 

- 1. Represent each document by vector of words
  - one attribute per word position in document
- 2. Learning: Use training examples to estimate
  - $\bullet P(+)$
  - $\bullet P(-)$
  - $\bullet P(doc|+)$
  - $\bullet P(doc|-)$

Naive Bayes conditional independence assumption

$$P(doc|v_j) = \prod_{i=1}^{length(doc)} P(a_i = w_k|v_j) \qquad \text{ith location}$$
 where  $P(a_i = w_k|v_j)$  is probability that word in position  $i$  is  $w_k$ , given  $v_j$  one more assumption:

one more assumption: 
$$P(a_i = w_k | v_j) = P(a_m = w_k | v_j), \forall i, m$$

#### Twenty NewsGroups

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

comp.graphics comp.os.ms-windows.misc comp.sys.ibm.pc.hardware comp.sys.mac.hardware comp.windows.x

misc.forsale rec.autos rec.motorcycles rec.sport.baseball rec.sport.hockey

alt.atheism
soc.religion.christian
talk.religion.misc
talk.politics.mideast
talk.politics.misc
talk.politics.misc

sci.space sci.crypt sci.electronics sci.med

Naive Bayes: 89% classification accuracy

#### Learn\_naive\_bayes\_text(Examples, V)

- 1. collect all words and other tokens that occur in Examples ~~50×
- $Vocabulary \leftarrow$  all distinct words and other tokens in Examples
  - 2. calculate the required  $P(v_j)$  and  $P(w_k|v_j)$  probability terms
  - For each target value  $v_j$  in V do
    - $-docs_j \leftarrow \text{subset of } Examples \text{ for which the }$ target value is  $v_j$
    - $-P(v_j) \leftarrow \frac{|docs_j|}{|Examples|}$
    - $-Text_j \leftarrow a \text{ single document created by } concatenating all members of <math>docs_j$

For code and data, see www.cs.cmu.edu/~tom/mlbook.html click on "Software and Data"

- $-n \leftarrow \text{total number of words in } Text_j \text{ (counting duplicate words multiple times)}$
- for each word  $w_k$  in Vocabulary
  - \*  $n_k \leftarrow$  number of times word  $w_k$  occurs in
  - $*P(w_k|v_j) \leftarrow \frac{n_k+1}{n+|Vocabulary|}$

#### CLASSIFY\_NAIVE\_BAYES\_TEXT(Doc)

- $positions \leftarrow$  all word positions in Doc that contain tokens found in Vocabulary
- Return  $v_{NB}$ , where

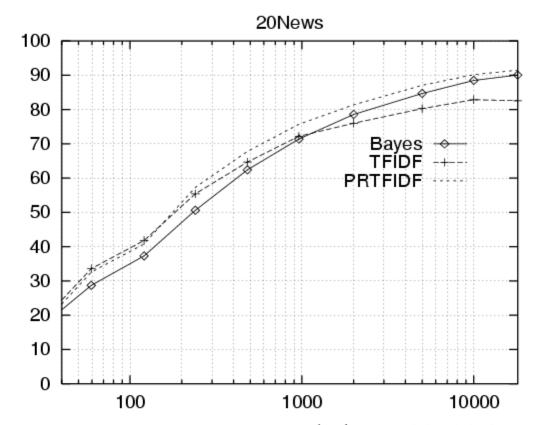
$$v_{NB} = \underset{v_{j} \in V}{\operatorname{argmax}} P(v_{j}) \prod_{i \in positions} P(a_{i}|v_{j})$$

$$V_{NB} = \underset{v_{j} \in V}{\operatorname{argmax}} (los(\mathcal{N})) \quad los P(v_{j}) + \underbrace{\geq los P(a_{i}|v_{j})}_{i}$$

$$P(V = v_{i}|Doc) = \underbrace{P(v_{i})}_{i} TP(a_{i}|v_{j})$$

$$P(Doc)$$

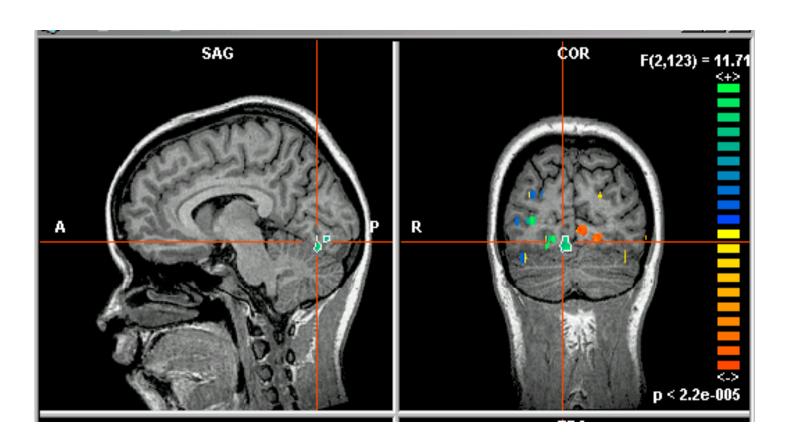
### Learning Curve for 20 Newsgroups



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

## What if we have continuous $X_i$ ?

Eg., image classification:  $X_i$  is real-valued ith pixel



## What if we have continuous $X_i$ ?

Eg., image classification:  $X_i$  is real-valued ith pixel

Naïve Bayes requires  $P(X_i | Y_k)$ , but  $X_i$  is real (continuous)

$$P(Y/x) = \frac{P(Y) P(x|Y)}{P(Y)} C.I. P(Y) T(P(x;|Y))$$

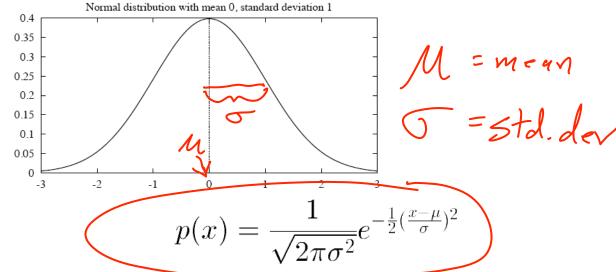
$$P(Y)$$

Common approach: assume  $P(X_i | Y_k)$  follows a normal (Gaussian) distribution

# Gaussian Distribution

(also known as "Normal" distribution)

p(x) is a probability density function, whose integral (not sum) is 1



The probability that X will fall into the interval (a, b) is given by

$$P(a(x)) = \int_a^b p(x) dx$$

• Expected, or mean value of X, E[X], is

$$E[X] = \mu$$

• Variance of X is

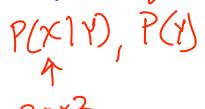
$$Var(X) = \sigma^2$$

• Standard deviation of X,  $\sigma_X$ , is

$$\sigma_X = \sigma$$

# What if we have continuous $X_i$ ?

Eg., image classification:  $X_i$  is ith pixel



Gaussian Naïve Bayes (GNB): assume

$$p(X_i = x | Y = y_k) = \frac{1}{\sqrt{2\pi\sigma_{ik}^2}} e^{-\frac{1}{2}(\frac{x-\mu_{ik}}{\sigma_{ik}})^2}$$
for single val of Y, 2n params for P(x)y) Mow many params?
Sometimes assume variance for Y boolean X=X,...X

- is independent of Y (i.e.,  $\sigma_i$ ),
- or independent of  $X_i$  (i.e.,  $\sigma_k$ ) • or both (i.e.,  $\sigma$ )

### Gaussian Naïve Bayes Algorithm – continuous X<sub>i</sub> (but still discrete Y)

 Train Naïve Bayes (examples) for each value  $y_k$ 

estimate\* 
$$\pi_k \equiv P(Y = y_k)$$

for each attribute  $X_i$  estimate class conditional mean  $\mu_{ik}$ , variance  $\sigma_{ik}$ 

• Classify  $(X^{new})$ 

$$Y^{new} \leftarrow \arg\max_{y_k} \ P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$
$$Y^{new} \leftarrow \arg\max_{y_k} \ \pi_k \prod_i Normal(X_i^{new}, \mu_{ik}, \sigma_{ik})$$

probabilities must sum to 1, so need estimate only n-1 parameters...

## Estimating Parameters: Y discrete, $X_i$ continuous

Maximum likelihood estimates:

jth training example

$$\widehat{\mu}_{ik} = \frac{1}{\sum_{j} \delta(Y^j = y_k)} \sum_{j} X_i^j \widehat{\delta(Y^j = y_k)}$$
 ith feature kth class

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

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

How many parameters must we estimate for Gaussian Naïve Bayes if Y has k possible values, X=<X1, ... Xn>?

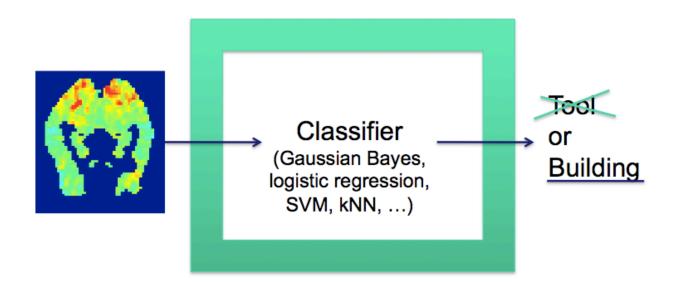
# What is form of decision surface for Gaussian Naïve Bayes classifier?

eg., if distributions are spherical, attributes have same variance (  $\sigma_{ik} = \sigma$  )

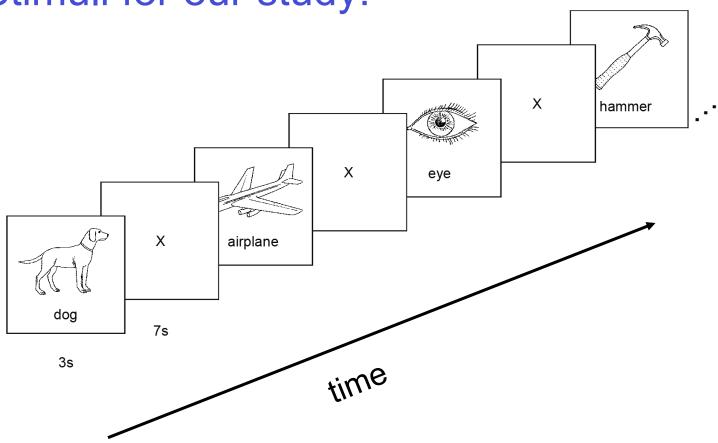
# What is form of decision surface for Naïve Bayes classifier?

# GNB Example: Classify a person's cognitive activity, based on brain image

- reading a sentence or viewing a picture?
- reading the word "Hammer" or "Apartment"
- viewing a vertical or horizontal line?
- answering the question, or getting confused?

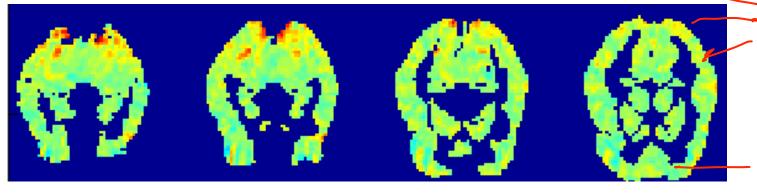


Stimuli for our study:

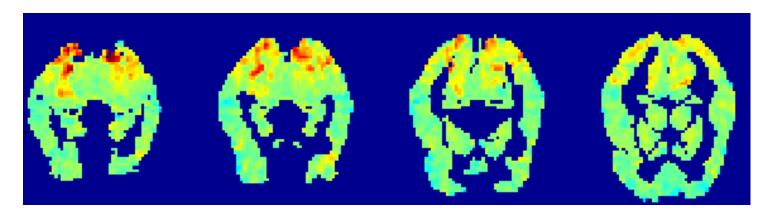


60 distinct exemplars, presented 6 times each

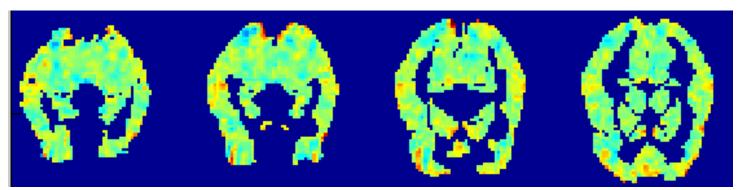
fMRI voxel means for "bottle": means defining P(Xi | Y="bottle)



Mean fMRI activation over all stimuli:



"bottle" minus mean activation:



high

robottle

**fMRI** 

activation

average

below average

### Training Classifiers over fMRI sequences

Learn the classifier function

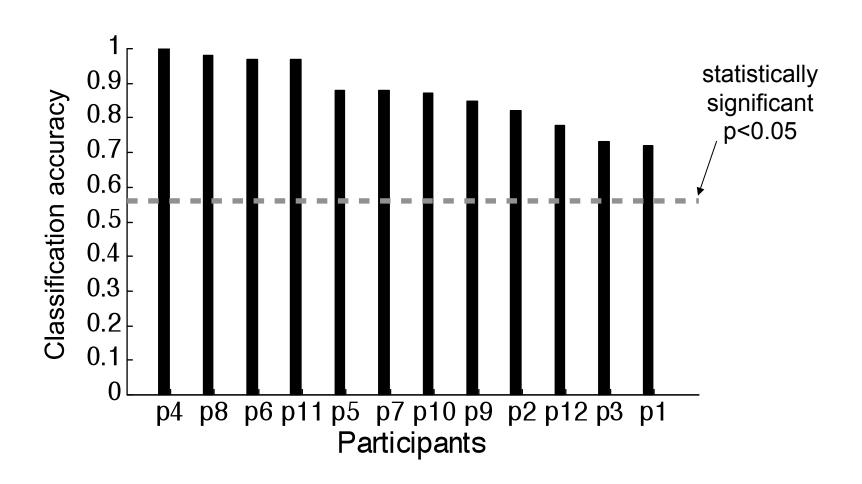
Mean(fMRI(t+4), ...,fMRI(t+7))  $\rightarrow$  WordCategory

Leave one out cross validation



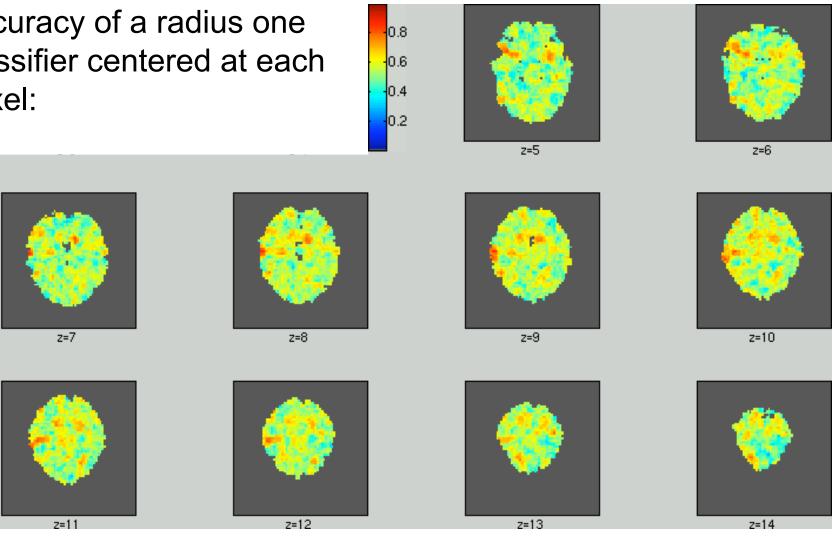
- Preprocessing:
  - Adjust for head motion
  - Adjust for nead motion Convert each image x to standard normal image  $x(i) \leftarrow \frac{x(i) \mu_x}{\sigma_x}$
- Learning algorithms tried:
  - kNN (spatial correlation)
  - SVM
  - SVDM
  - Gaussian Naïve Bayes
  - Regularized Logistic regression
- Feature selection methods tried:
  - Logistic regression weights, voxel stability, activity relative to fixation,...

### Classification task: is person viewing a "tool" or "building"?



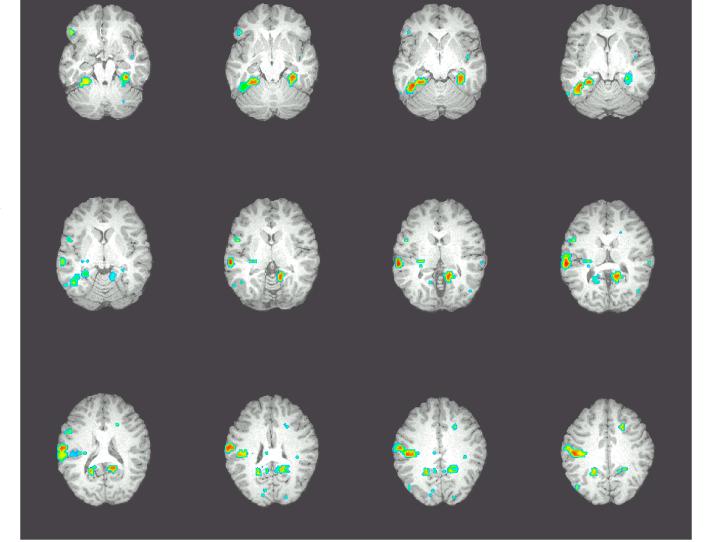
## Where in the brain is activity that distinguishes tools vs. buildings?

Accuracy of a radius one classifier centered at each voxel:



## voxel clusters: searchlights

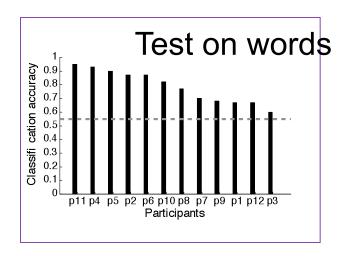
Accuracies of cubical 27-voxel classifiers centered at each significant voxel [0.7-0.8]

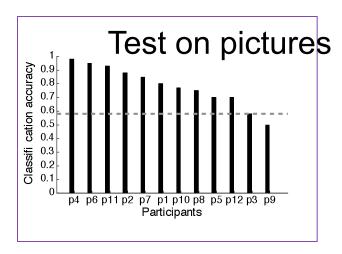


# Are classifiers detecting neural representations of meaning or perceptual features?

ML: Can we train on word stimuli, then decode picture stimuli?

YES: We can train classifiers when presenting English words, then decode category of picture stimuli, or Portuguese words



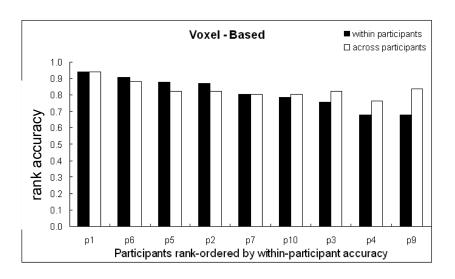


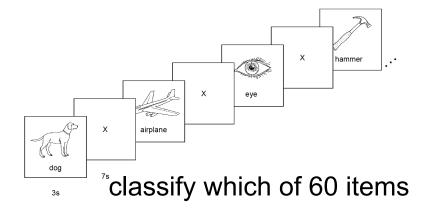
Therefore, the learned neural activation patterns must capture how the brain represents the <u>meaning</u> of input stimulus

### Are representations similar across different people?

ML: Can we train classifier on data from a collection of people, then decode stimuli for a new person?

YES: We can train on one group of people, and classify fMRI images of new person





Therefore, seek a theory of neural representations common to all of us (and of how we vary)

## What you should know:

- Training and using classifiers based on Bayes rule
- Conditional independence
  - What it is
  - Why it's important
- Naïve Bayes
  - What it is
  - Why we use it so much
  - Training using MLE, MAP estimates
  - Discrete variables (Bernoulli) and continuous (Gaussian)

## Questions to think about:

 Can you use Naïve Bayes for a combination of discrete and real-valued X<sub>i</sub>?

 How can we easily model just 2 of n attributes as dependent?

 What does the decision surface of a Naïve Bayes classifier look like?

How would you select a subset of X<sub>i</sub>'s?