# Machine Learning

KMA Solaiman
UMBC CMSC 471

Today:
- Naïve Bayes
    - discrete-valued $X_i$'s
    - Document classification
- Gaussian Naïve Bayes
    - real-valued $X_i$'s
    - Brain image classification

Recently:

- Bayes classifiers to learn P(Y|X)
- MLE and MAP estimates for parameters of P
- Conditional independence
- Naïve Bayes → make Bayesian learning practical

Next:

- Text classification
- Naïve Bayes and continuous variables $X_i$:
  - Gaussian Naïve Bayes classifier
- Learn P(Y|X) directly
  - Logistic regression, Regularization, Gradient ascent
- Naïve Bayes or Logistic Regression?
  - Generative vs. Discriminative classifiers

# Estimating Parameters

- Maximum Likelihood Estimate (MLE): choose θ that maximizes probability of observed data $\mathcal{D}$

$$\widehat{\theta} = \arg\max_{\theta} \; P(\mathcal{D} \mid \theta)$$

- Maximum a Posteriori (MAP) estimate: choose θ that is most probable given prior probability and the data

$$\widehat{\theta} = \arg\max_{\theta} \; P(\theta \mid \mathcal{D})$$

$$= \arg\max_{\theta} \; = \; \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

# Estimating Parameters: $Y, X_i$ discrete-valued

Maximum likelihood estimates:

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \wedge Y = y_k\}}{\#D\{Y = y_k\}}$$

MAP estimates (Beta, Dirichlet priors):

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\} + (\beta_k - 1)}{|D| + \sum_m (\beta_m - 1)}$$

Only difference: "imaginary" examples

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \wedge Y = y_k\} + (\beta_k - 1)}{\#D\{Y = y_k\} + \sum_m (\beta_m - 1)}$$

# Naïve Bayes: classifying text documents

- Classify which emails are spam?
- Classify which emails promise an attachment?

I am pleased to announce that Bob Frederking of the Language Technologies Institute is our new Associate Dean for Graduate Programs. In this role, he oversees the many issues that arise with our multiple masters and PhD programs. Bob brings to this position considerable experience with the masters and PhD programs in the LTI.

I would like to thank Frank Pfenning, who has served ably in this role for the past two years.

**********************************
Randal E. Bryant
Dean and University Professor

How shall we represent text documents for Naïve Bayes?

# Learning to classify documents: P(Y|X)

- Y discrete valued.
  - e.g., Spam or not
- X = <$X_1$, $X_2$, … $X_n$> = document

> I am pleased to announce that Bob Frederking of the Language Technologies Institute is our new Associate Dean for Graduate Programs. In this role, he oversees the many issues that arise with our multiple masters and PhD programs. Bob brings to this position considerable experience with the masters and PhD programs in the LTI.
>
> I would like to thank Frank Pfenning, who has served ably in this role for the past two years.
>
> ********************************
> Randal E. Bryant
> Dean and University Professor

- $X_i$ is a random variable describing…

# Learning to classify documents: P(Y|X)

- Y discrete valued.
  - e.g., Spam or not
- X = $<X_1, X_2, \ldots X_n>$ = document

I am pleased to announce that Bob Frederking of the Language Technologies Institute is our new Associate Dean for Graduate Programs. In this role, he oversees the many issues that arise with our multiple masters and PhD programs. Bob brings to this position considerable experience with the masters and PhD programs in the LTI.

I would like to thank Frank Pfenning, who has served ably in this role for the past two years.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

Randal E. Bryant
Dean and University Professor

- $X_i$ is a random variable describing…

Answer 1: $X_i$ is boolean, 1 if word i is in document, else 0

$$\text{e.g., } X_{pleased} = 1$$

Issues?

# Learning to classify documents: P(Y|X)

- Y discrete valued.
  - e.g., Spam or not

- X = $<X_1, X_2, \ldots X_n>$ = document

> I am pleased to announce that Bob Frederking of the Language Technologies Institute is our new Associate Dean for Graduate Programs. In this role, he oversees the many issues that arise with our multiple masters and PhD programs. Bob brings to this position considerable experience with the masters and PhD programs in the LTI.
>
> I would like to thank Frank Pfenning, who has served ably in this role for the past two years.
>
> \*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*
>
> Randal E. Bryant
> Dean and University Professor

- $X_i$ is a random variable describing…

Answer 2:

- $X_i$ represents the $i^{th}$ *word position* in document
- $X_1$ = "I", $X_2$ = "am", $X_3$ = "pleased"
- and, let's assume the $X_i$ are iid (indep, identically distributed)

$$P(X_i|Y) = P(X_j|Y) \quad (\forall i, j)$$

# Learning to classify document: P(Y|X) the "Bag of Words" model

- Y discrete valued.  e.g., Spam or not

- $X = <X_1, X_2, \ldots X_n>$ = document

- $X_i$ are iid random variables.  Each represents the word at its position i in the document

- Generating a document according to this distribution = rolling a 50,000 sided die, once for each word position in the document

- The observed counts for each word follow a ??? distribution

# Multinomial Distribution

Eg. 2  Dice roll problem (6 outcomes instead of 2)

Likelihood is ~ Multinomial($\theta = \{\theta_1, \theta_2, \ldots, \theta_k\}$)

$$P(\mathcal{D} \mid \theta) = \theta_1^{\alpha_1}\theta_2^{\alpha_2}\ldots\theta_k^{\alpha_k}$$

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\prod_{i=1}^{k}\theta_i^{\beta_i-1}}{B(\beta_1,\ldots,\beta_k)} \sim \text{Dirichlet}(\beta_1,\ldots,\beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta|D) \sim \text{Dirichlet}(\beta_1 + \alpha_1,\ldots,\beta_k + \alpha_k)$$

$$\hat{\theta_i}^{MAP} = \hat{P}(X = i) = \frac{\alpha_i + \beta_i - 1}{\sum_{j=1}^{k}(\alpha_j + \beta_j - 1)}$$

# Multinomial Distribution

K sides ↓

**Eg. 2** Dice roll problem (6 outcomes instead of 2)

Likelihood is ~ Multinomial($\theta = \{\theta_1, \theta_2, \ldots, \theta_k\}$)

$$P(\mathcal{D} \mid \theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \ldots \theta_k^{\alpha_k}$$

$\theta_i = P(X = i)$

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\prod_{i=1}^{k} \theta_i^{\beta_i - 1}}{B(\beta_1, \ldots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \ldots, \beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta|D) \sim \text{Dirichlet}(\beta_1 + \alpha_1, \ldots, \beta_k + \alpha_k)$$

$$\hat{\theta_i}^{MAP} = \hat{P}(X = i) = \frac{\alpha_i + \beta_i - 1}{\sum_{j=1}^{k}(\alpha_j + \beta_j - 1)}$$

*Dirichlet is a generalization of Beta Distribution*

# Multinomial Bag of Words

counts
$\alpha_i$'s



| the world of ► All About The Company | |
|---|---|
| **TOTAL** | Global Activities |
| | Corporate Structure |
| | TOTAL's Story |
| | Upstream Strategy |
| | Downstream Strategy |
| | Chemicals Strategy |
| | TOTAL Foundation |
| | Homepage |

all about the
**company**

Our energy exploration, production, and distribution operations span the globe, with activities in more than 100 countries.

At TOTAL, we draw our greatest strength from our fast-growing oil and gas reserves. Our strategic emphasis on natural gas provides a strong position in a rapidly expanding market.

Our expanding refining and marketing operations in Asia and the Mediterranean Rim complement already solid positions in Europe, Africa, and the U.S.

Our growing specialty chemicals sector adds balance and profit to the core energy business.

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

Bag of Word Model

# MAP estimates for bag of words

Map estimate for multinomial

$$\hat{\theta}_i^{MAP} = \hat{P}(X = i) = \frac{\alpha_i + \beta_i - 1}{\sum_{j=1}^{k}(\alpha_j + \beta_j - 1)}$$

$$\hat{\theta}_{aardvark}^{MAP} = P(X = \text{aardvark}) = \frac{\#\text{ observed 'aardvark' } + \#\text{ hallucinated 'aardvark'}}{\#\text{ observed words } + \#\text{ hallucinated words}}$$

What $\beta$'s should we choose?

# MAP estimates for bag of words

Map estimate for multinomial

$$\hat{\theta}_i^{MAP} = \hat{P}(X = i) = \frac{\alpha_i + \beta_i - 1}{\sum_{j=1}^{k}(\alpha_j + \beta_j - 1)}$$

$$\hat{\theta}_{aardvark}^{MAP} = P(X = \text{aardvark}) = \frac{\# \text{ observed 'aardvark'} + \# \text{ hallucinated 'aardvark'}}{\# \text{ observed words} + \# \text{ hallucinated words}}$$

What $\beta$'s should we choose?

- Large document, how many times each word occur
- Uniform distribution

# Naïve Bayes Algorithm – discrete $X_i$

- Train Naïve Bayes (examples)

  for each value $y_k$

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

      for each value $x_{ij}$ of each attribute $X_i$

          estimate $\theta_{ijk} \equiv P(X_i = x_{ij} | Y = y_k)$

  prob that word $x_{ij}$ appears in position i, given Y=$y_k$

- Classify ($X^{new}$)

$$Y^{new} \leftarrow \arg\max_{y_k} \; P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg\max_{y_k} \; \pi_k \prod_i \theta_{ijk}$$

[*] Additional assumption: word probabilities are position independent

$$\theta_{ijk} = \theta_{mjk} \;\; \text{for} \;\; i \neq m$$

# Naïve Bayes Algorithm – discrete $X_i$

- Train Naïve Bayes (examples)

  for each value $y_k$

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

      Spam: $k = 1$

      $\neg$ spam: $k = 0$

      for each value $x_{ij}$ of each attribute $X_i$

          estimate $\theta_{ijk} \equiv P(X_i = x_{ij} | Y = y_k)$

  <div style="background-color:#c4c9f0">prob that word x<sub>ij</sub> appears in position i, given Y=y<sub>k</sub></div>

- Classify ($X^{new}$)

$$Y^{new} \leftarrow \arg\max_{y_k} \; P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg\max_{y_k} \; \pi_k \prod_i \theta_{ijk}$$

*  Additional assumption:  word probabilities are position independent

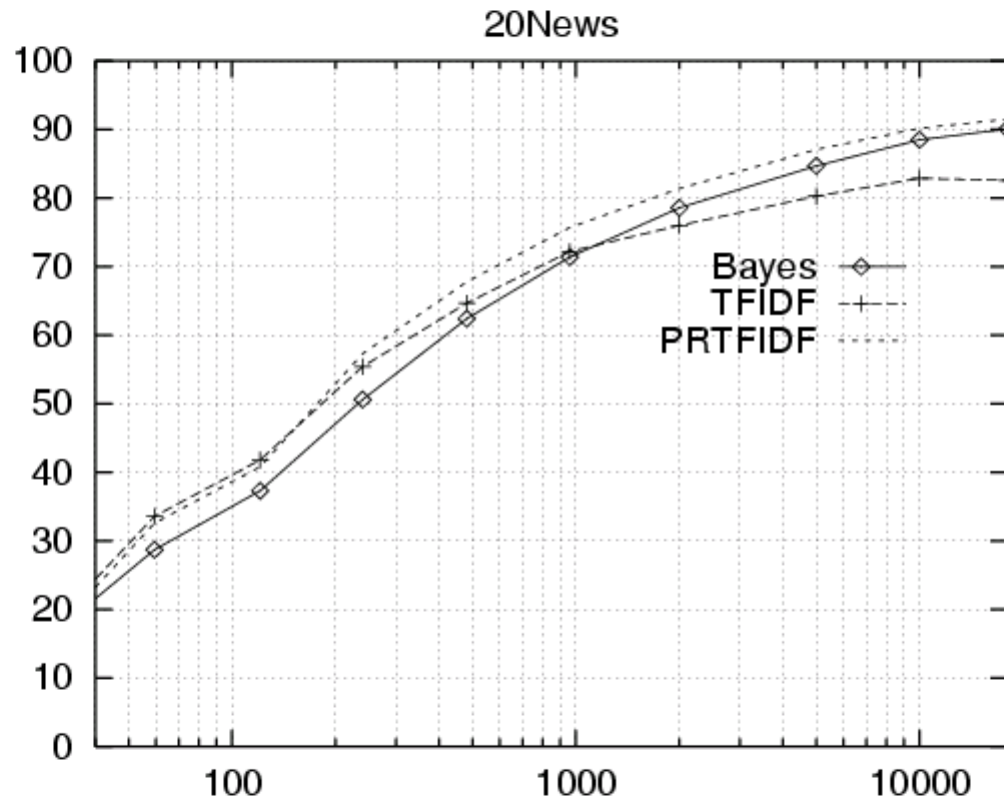$$\theta_{ijk} = \theta_{mjk} \quad \text{for} \quad i \neq m$$

# Twenty NewsGroups

---

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

# Learning Curve for 20 Newsgroups



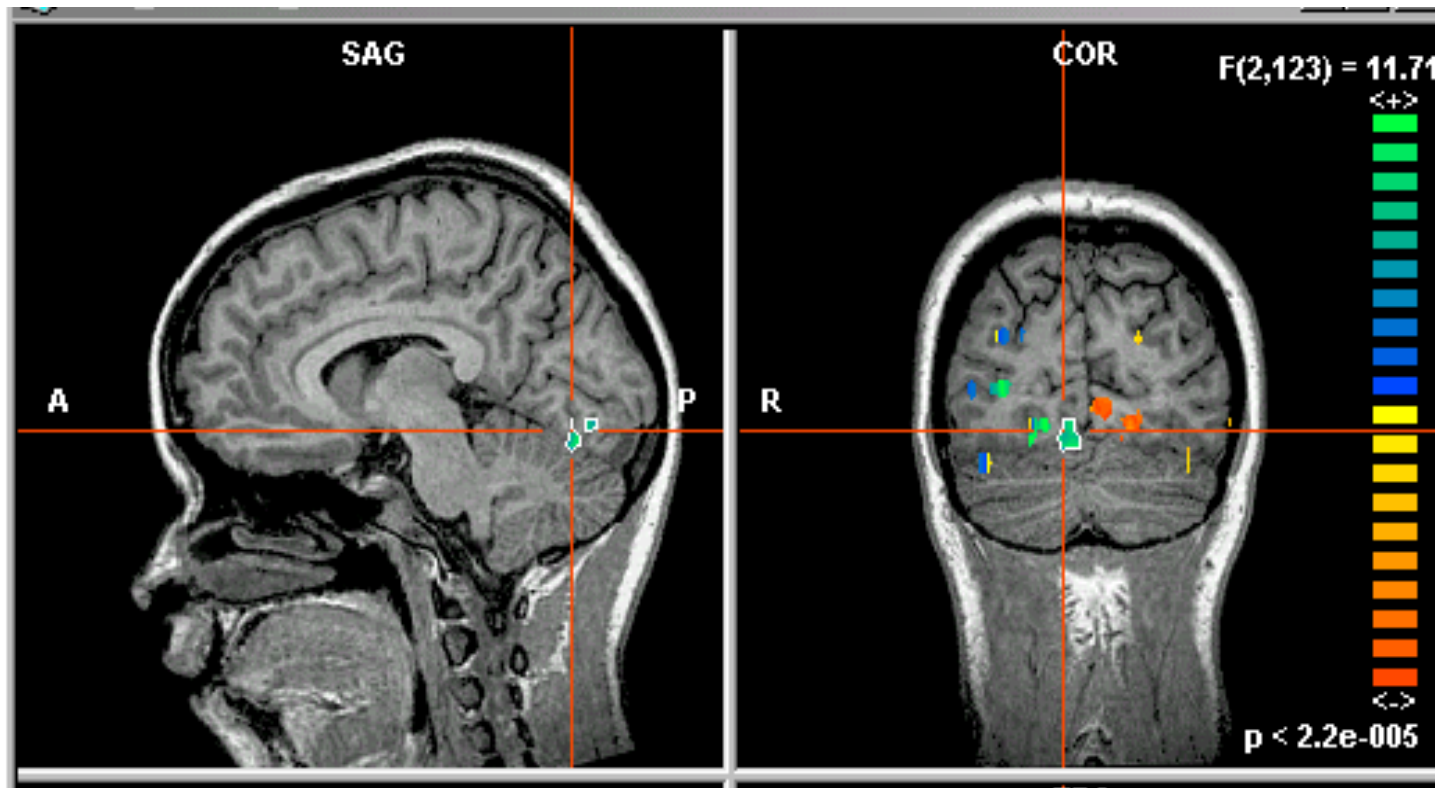Accuracy vs. Training set size (1/3 withheld for test)

# 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 i<sup>th</sup> pixel

# What if we have continuous $X_i$ ?

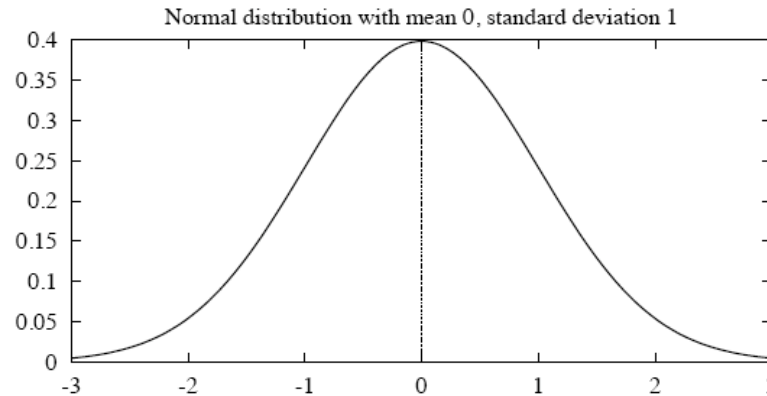Eg., image classification: $X_i$ is real-valued $i^{th}$ pixel

Naïve Bayes requires $P(X_i \mid Y=y_k)$, but $X_i$ is real (continuous)

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

Common approach: assume $P(X_i \mid Y=y_k)$ follows a Normal (Gaussian) distribution

# What if we have continuous $X_i$?

Eg., image classification: $X_i$ is real-valued i[th] pixel

Naïve Bayes requires $P(X_i | Y=y_k)$, but $X_i$ is real (continuous)

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

Common approach: assume $P(X_i | Y=y_k)$ follows a Normal (Gaussian) distribution

# Gaussian Distribution
(also called "Normal")

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

Normal distribution with mean 0, standard deviation 1

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

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

$$\int_a^b p(x)dx$$

- 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 $X_i$ ?

Gaussian Naïve Bayes (GNB): assume

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

Sometimes assume variance
- 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 $X_i$
## (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 $P(X_i | Y = y_k)$

  - class conditional mean $\mu_{ik}$, variance $\sigma_{ik}$

- Classify ($X^{new}$)

$$Y^{new} \leftarrow \arg\max_{y_k} \; P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg\max_{y_k} \; \pi_k \prod_i \mathcal{N}(X_i^{new}; \mu_{ik}, \sigma_{ik})$$

\* 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

$$\widehat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k)$$

ith feature

kth class
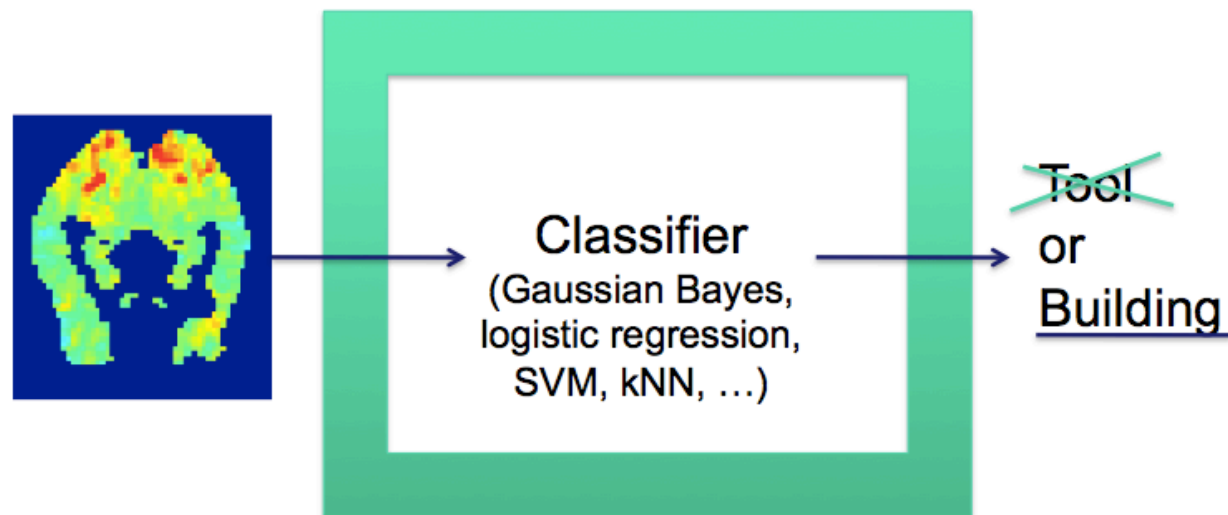
$\delta()=1$ if $(Y^j = y_k)$ else $0$

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

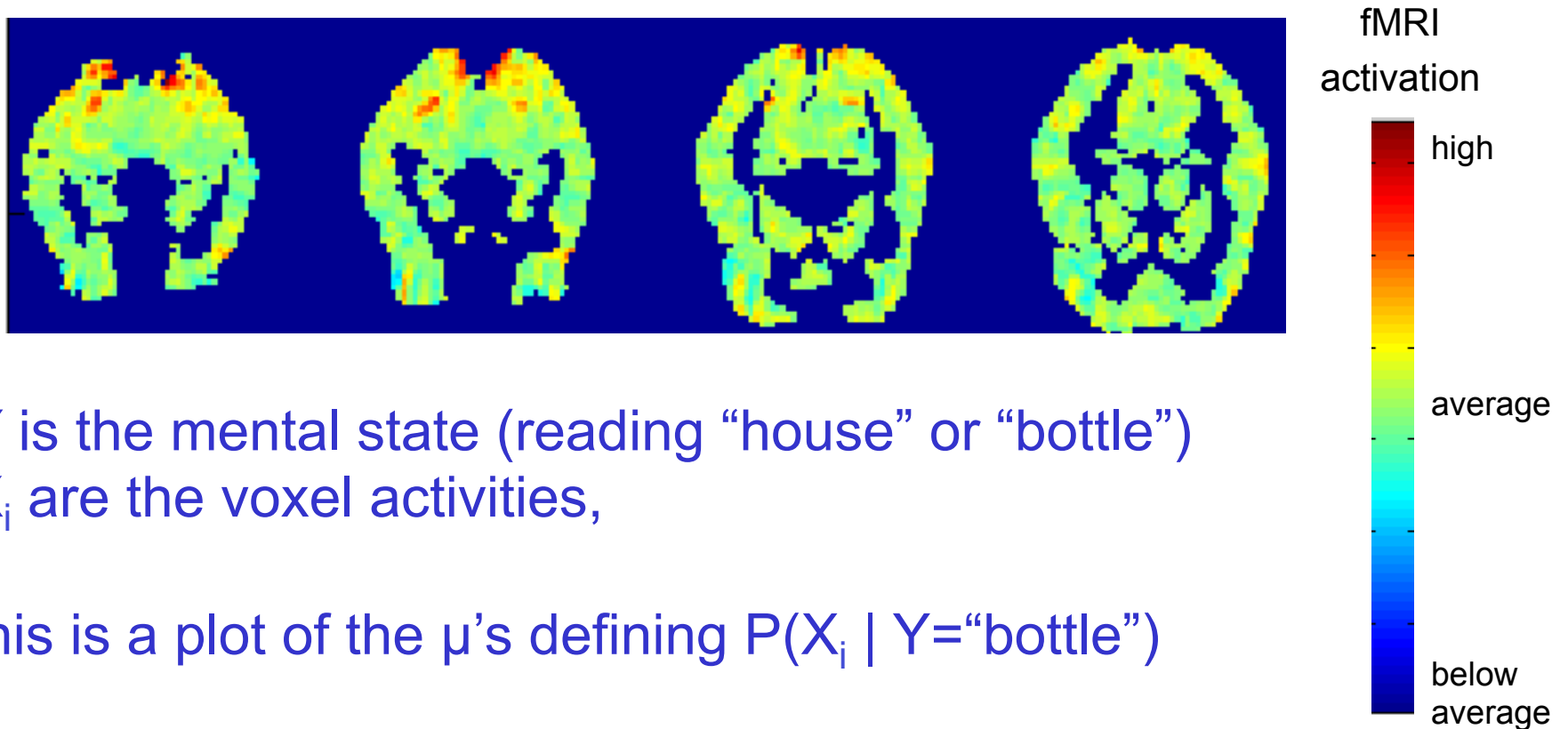How many parameters must we estimate for Gaussian Naïve Bayes if Y has k possible values, X=<X1, … Xn>?

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

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

- reading a sentence or viewing a picture?
- reading the word describing a "Tool" or "Building"?
- answering the question, or getting confused?

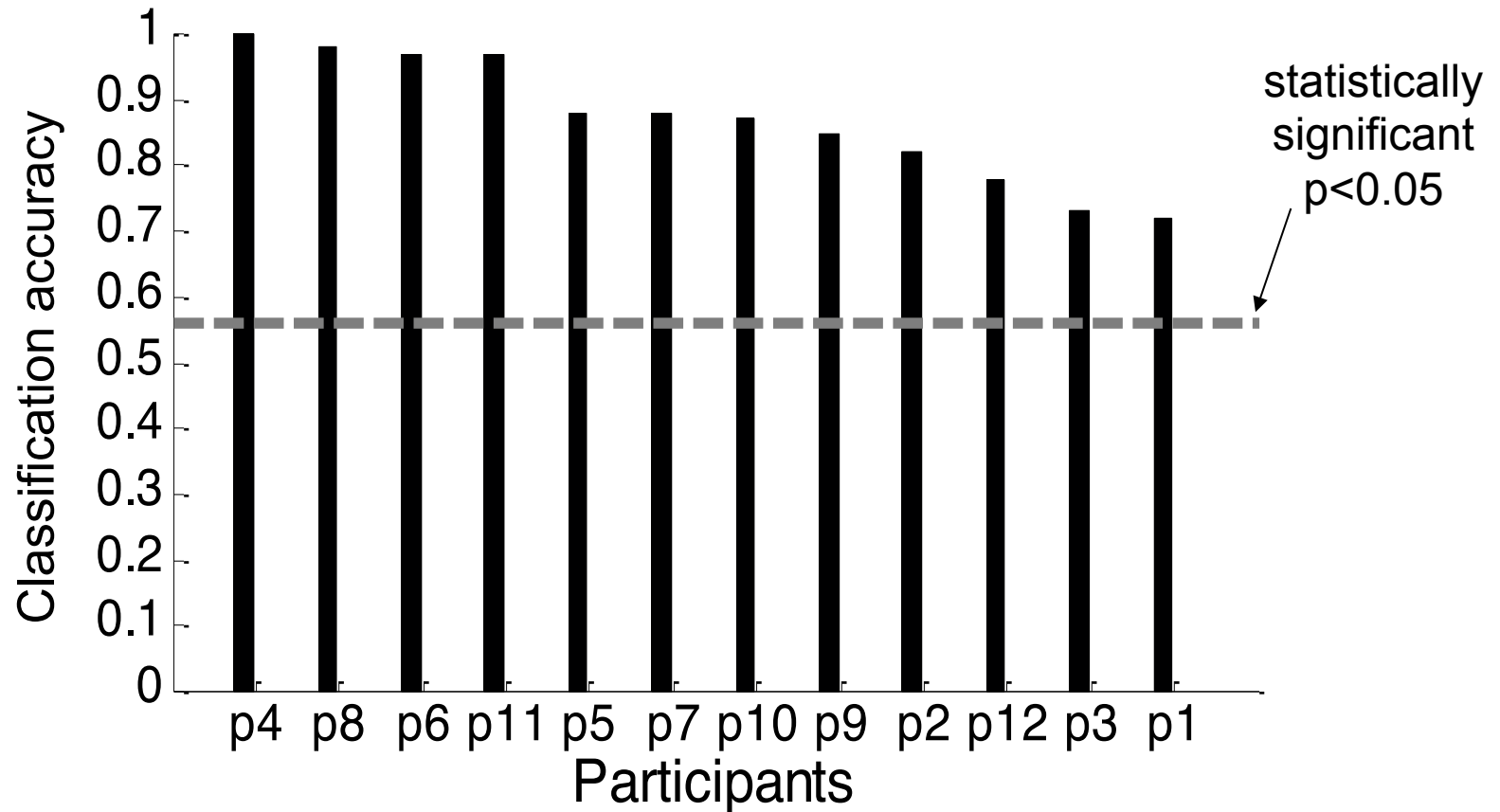# Mean activations over all training examples for Y="bottle"



Y is the mental state (reading "house" or "bottle")
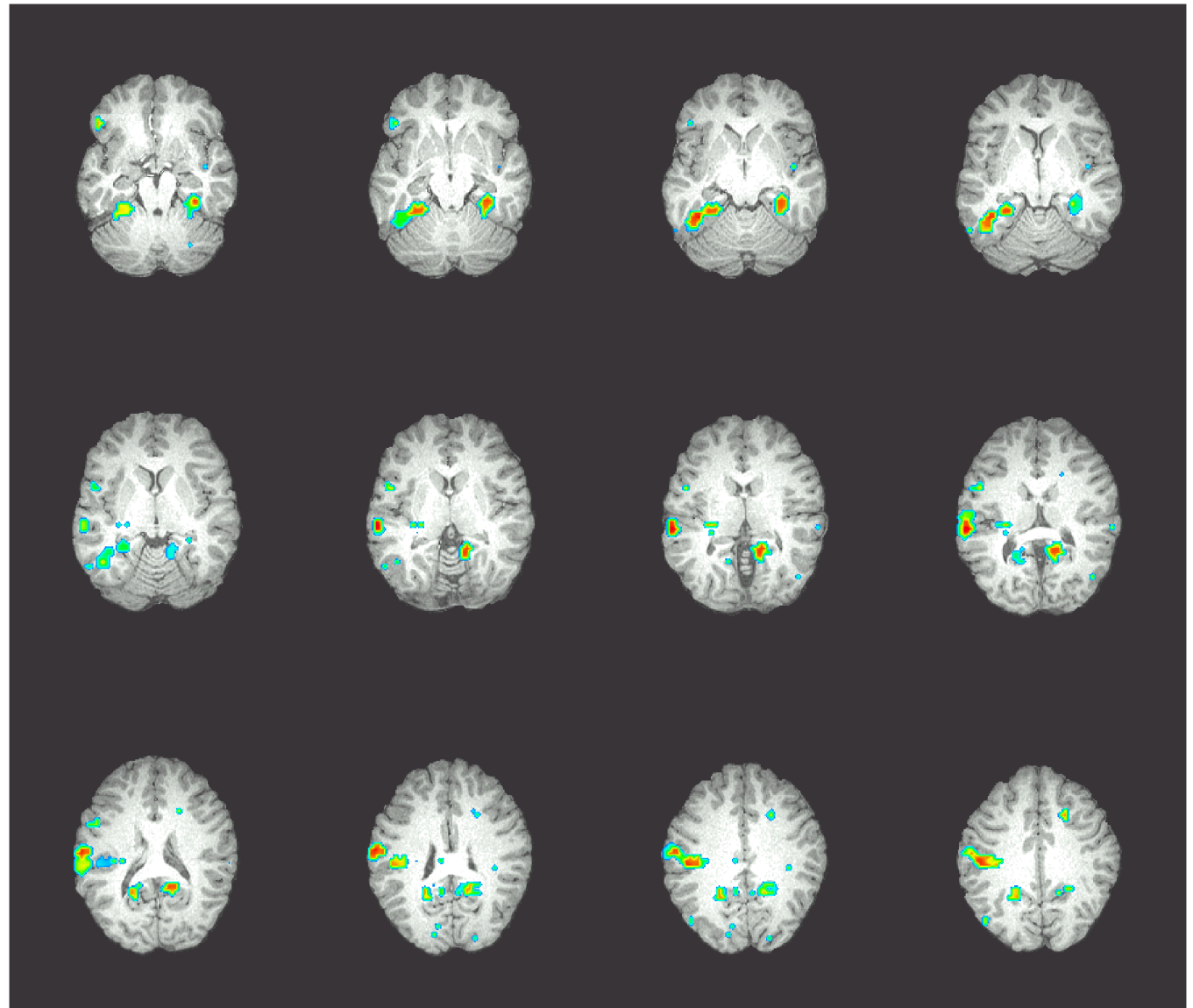$X_i$ are the voxel activities,

this is a plot of the $\mu$'s defining $P(X_i \mid Y=$"bottle"$)$

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

# Where is information encoded in the brain?

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

# Naïve Bayes: What you should know

- Designing classifiers based on Bayes rule

- Conditional independence
  - What it is
  - Why it's important

- Naïve Bayes assumption and its consequences
  - Which (and how many) parameters must be estimated under different generative models (different forms for P(X|Y) )
    - and why this matters

- How to train Naïve Bayes classifiers
  - MLE and MAP estimates
  - with discrete and/or continuous inputs $X_i$

# Questions to think about:

- Can you use Naïve Bayes for a combination of discrete and real-valued $X_i$?

- 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 $X_i$'s?