

# Machine Learning 10-601

Tom M. Mitchell  
Machine Learning Department  
Carnegie Mellon University

September 17, 2017

## Today:

- Gaussian Naïve Bayes
- Logistic regression
- Gradient ascent/descent

## Readings: (see class website)

### Required:

- Mitchell: “Naïve Bayes and Logistic Regression”

# What if we have continuous $X_i$ ?

Eg., image classification:  $X_i$  is real-valued i<sup>th</sup> pixel

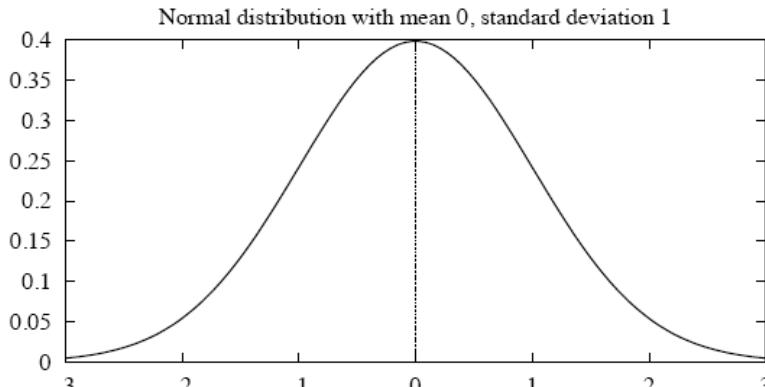
Naïve Bayes requires  $P(X_i | Y=y_k)$ , but  $X_i$  is real (continuous)

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

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



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

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

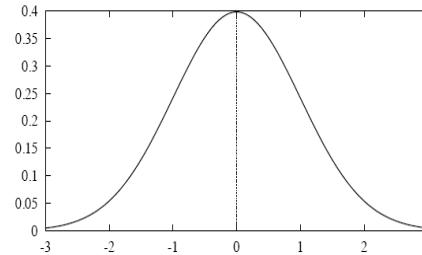
$$Var(X) = \sigma^2$$

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

$$\sigma_X = \sigma$$

# What if we have continuous $X_i$ ?

Gaussian Naïve Bayes (GNB): assume



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

Sometimes assume variance  $\sigma_{ik}$

- 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:

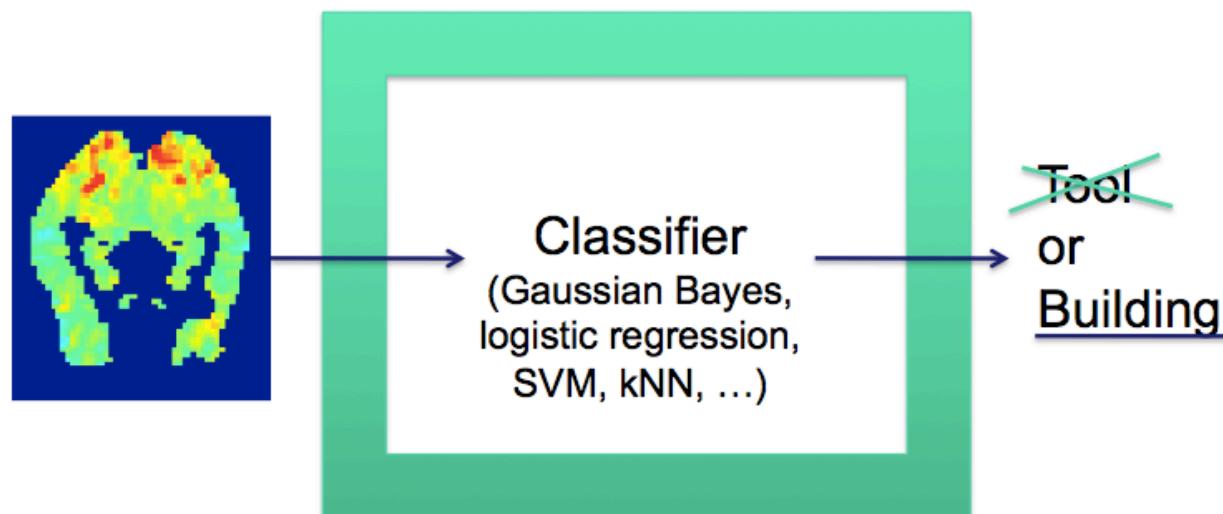
$$\hat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k)$$

ith feature      kth class      jth training example  
 $\delta()=1$  if  $(Y^j=y_k)$   
else 0

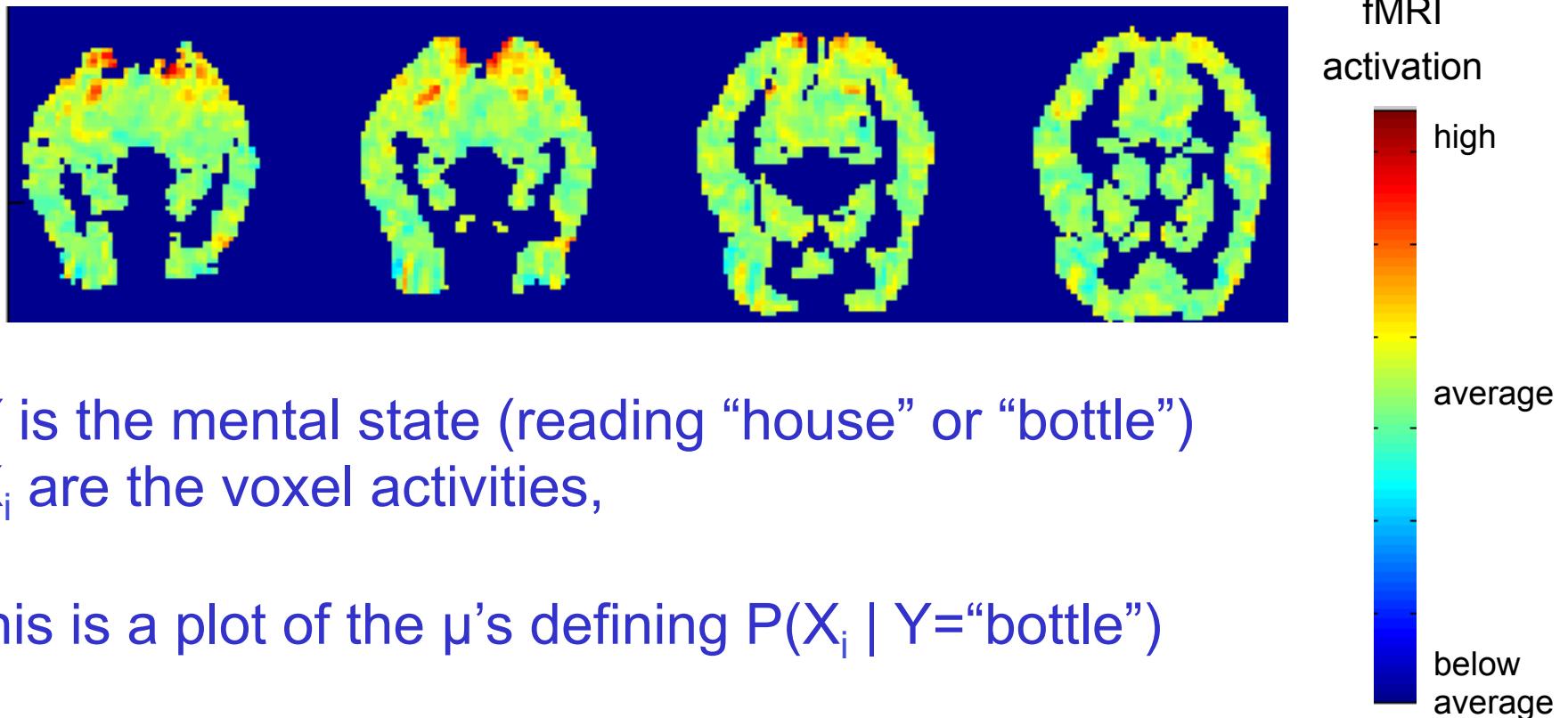
$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j (X_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k)$$

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

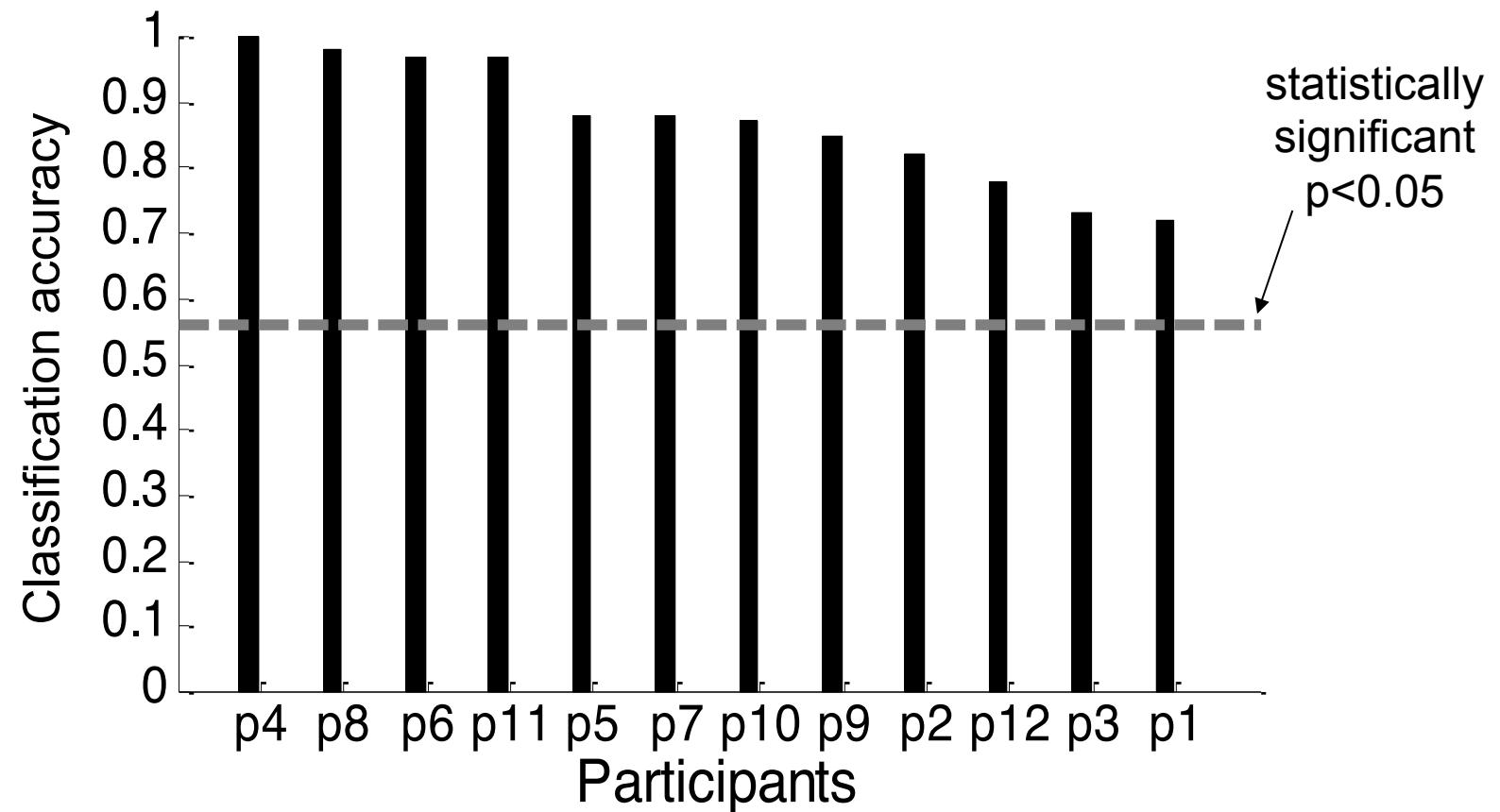
- reading a sentence or viewing a picture?
- reading the word describing a “Tool” or “Building”?
- answering the question, or getting confused?



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

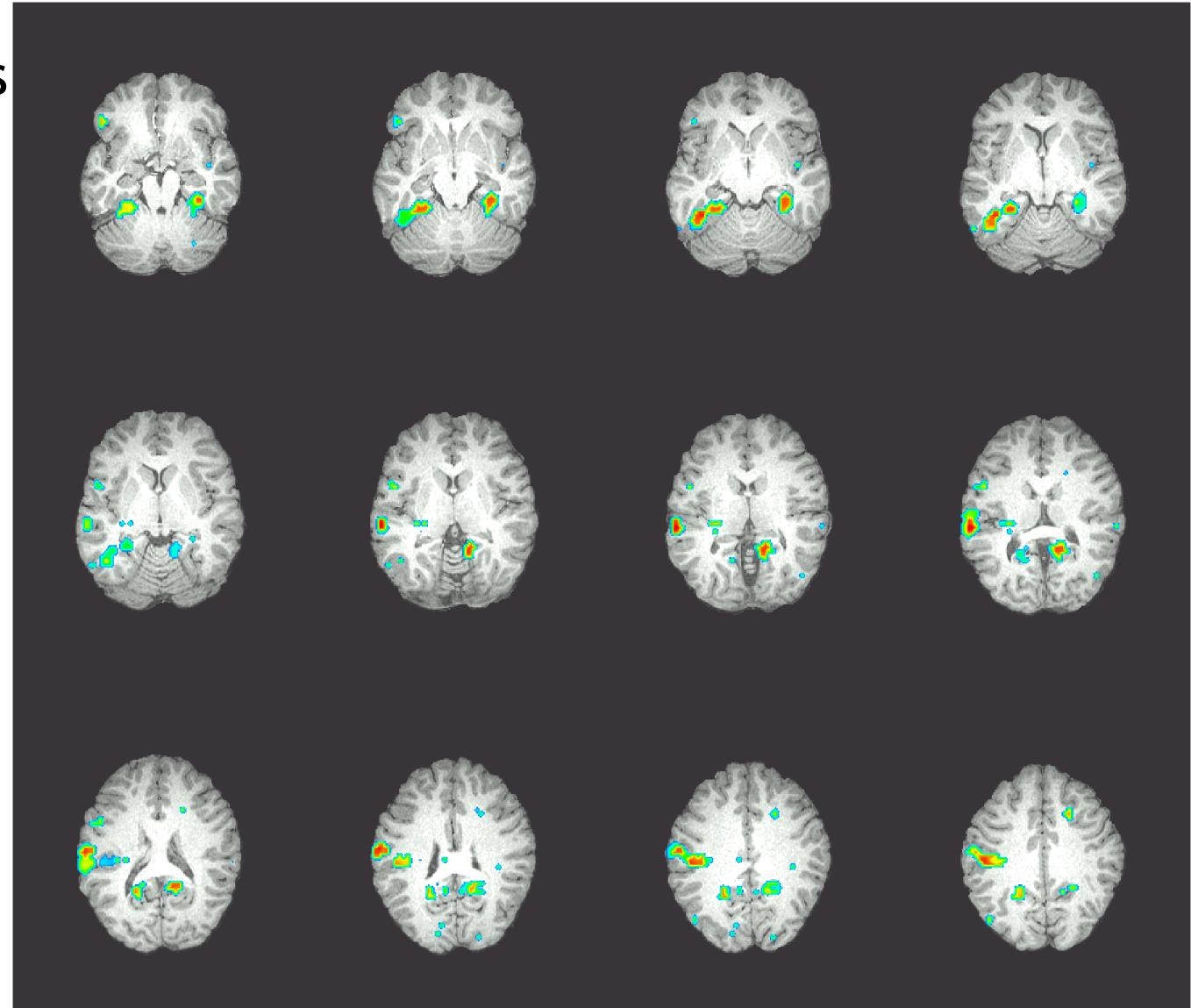


## Classification task: is person viewing a “tool” or “building”?



# Where is information encoded in the brain?

Tools vs. Buildings  
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 extend Naïve Bayes to model that two of the  $X_i$ 's are dependent?
- What does the decision surface of a Naïve Bayes classifier look like?
- What error will Naïve Bayes achieve if cond. indep. is satisfied and we have infinite training data?

# Logistic Regression

Idea:

- Naïve Bayes allows computing  $P(Y|X)$  by learning  $P(Y)$  and  $P(X|Y)$
- Why not learn  $P(Y|X)$  directly?

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

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

Derive form for  $P(Y|X)$  for Gaussian  $P(X_i|Y=y_k)$  assuming  $\sigma_{ik} = \sigma_i$

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

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

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

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

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

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

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

# 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

$$P(Y = 0|X = \langle X_1, \dots, X_n \rangle) =$$

implies

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

implies

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

# 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

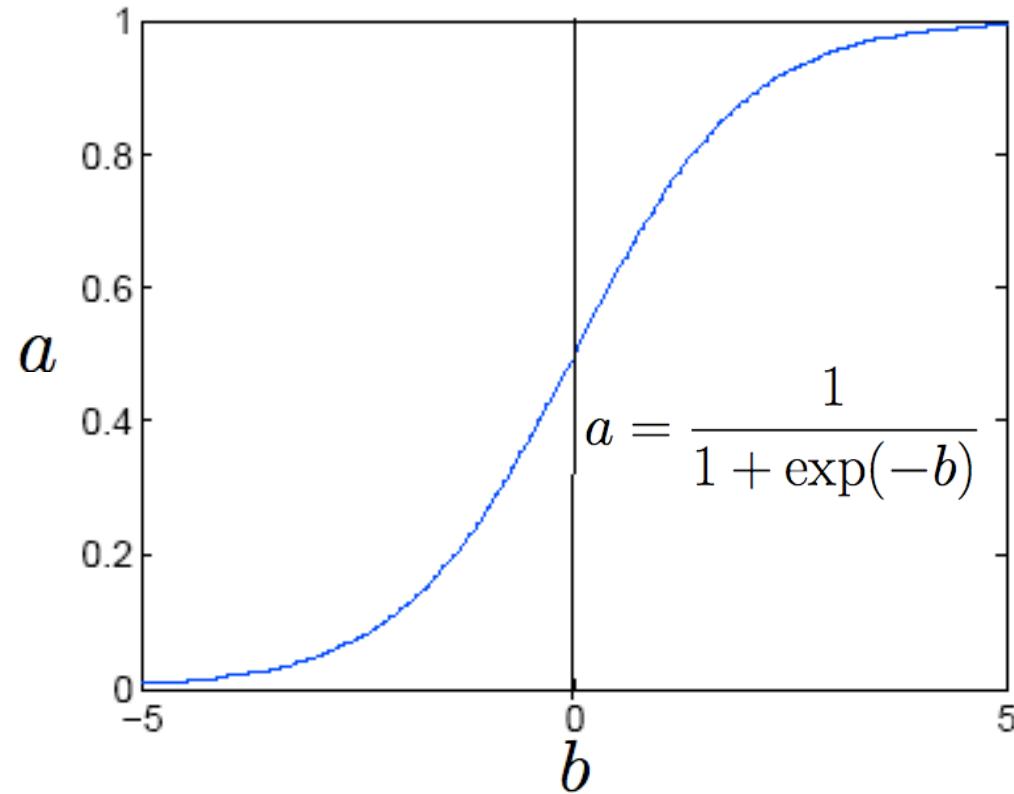
$$\frac{P(Y = 0|X)}{P(Y = 1|X)} = \exp(w_0 + \sum_i w_i X_i)$$

linear  
classification  
rule!

implies

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

# Logistic function



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

# Logistic regression more generally

- Logistic regression when Y not boolean (but still discrete-valued).
- Now  $y \in \{y_1 \dots y_R\}$  : 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  $k = R$      $P(Y = y_R | X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji}X_i)}$

# Training Logistic Regression: MLE, MCLE

- we have L training examples:  $\{\langle X^1, Y^1 \rangle, \dots \langle X^L, Y^L \rangle\}$

- maximum likelihood estimate for parameters W

$$\begin{aligned} W_{MLE} &= \arg \max_W P(\langle X^1, Y^1 \rangle \dots \langle X^L, Y^L \rangle | W) \\ &= \arg \max_W \prod_l P(\langle X^l, Y^l \rangle | W) \end{aligned}$$

- maximum conditional likelihood estimate...

# Training Logistic Regression: MCLE

- Choose parameters  $W = \langle w_0, \dots, w_n \rangle$  to maximize conditional likelihood of training data

where

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

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

- Training data  $D = \{\langle X^1, Y^1 \rangle, \dots, \langle X^L, Y^L \rangle\}$
- Data likelihood =  $\prod_l P(X^l, Y^l | W)$
- Data conditional likelihood =  $\prod_l P(Y^l | X^l, W)$

$$W_{MCLE} = \arg \max_W \prod_l P(Y^l | W, X^l)$$

# Expressing Conditional Log Likelihood

$$l(W) \equiv \ln \prod_l P(Y^l | X^l, W) = \sum_l \ln P(Y^l | X^l, W)$$

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

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

$$\begin{aligned} l(W) &= \sum_l Y^l \ln P(Y^l = 1 | X^l, W) + (1 - Y^l) \ln P(Y^l = 0 | X^l, W) \\ &= \sum_l Y^l \ln \frac{P(Y^l = 1 | X^l, W)}{P(Y^l = 0 | X^l, W)} + \ln P(Y^l = 0 | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l)) \end{aligned}$$

# Maximizing Conditional Log Likelihood

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

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

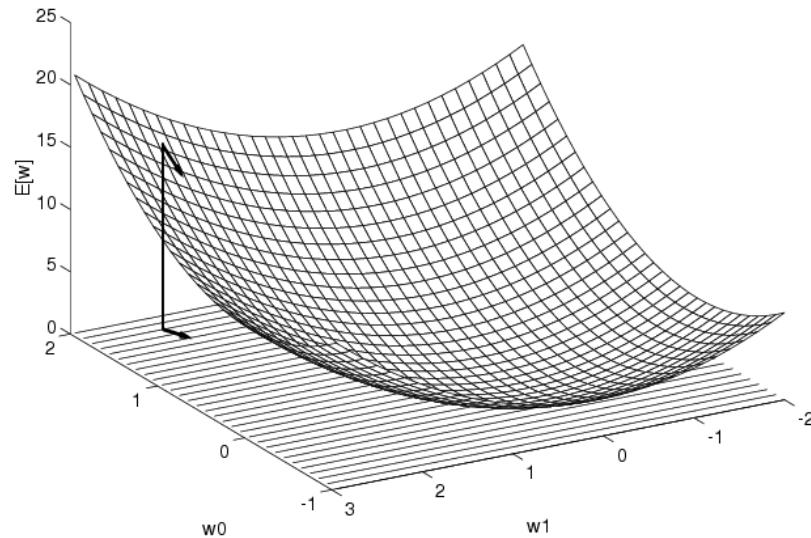
$$\begin{aligned} l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) = \sum_l \ln P(Y^l | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l)) \end{aligned}$$

Good news:  $l(W)$  is concave function of  $W$

Bad news: no closed-form solution to maximize  $l(W)$

# Gradient Descent

---



Gradient

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

i.e.,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

# Gradient Descent:

*Batch gradient:* use error  $E_D(\mathbf{w})$  over entire training set  $D$

Do until satisfied:

1. Compute the gradient  $\nabla E_D(\mathbf{w}) = \left[ \frac{\partial E_D(\mathbf{w})}{\partial w_0} \dots \frac{\partial E_D(\mathbf{w})}{\partial w_n} \right]$
2. Update the vector of parameters:  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla E_D(\mathbf{w})$

*Stochastic gradient:* use error  $E_d(\mathbf{w})$  over single examples  $d \in D$

Do until satisfied:

1. Choose (with replacement) a random training example  $d \in D$
2. Compute the gradient just for  $d$ :  $\nabla E_d(\mathbf{w}) = \left[ \frac{\partial E_d(\mathbf{w})}{\partial w_0} \dots \frac{\partial E_d(\mathbf{w})}{\partial w_n} \right]$
3. Update the vector of parameters:  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla E_d(\mathbf{w})$

Stochastic approximates Batch arbitrarily closely as  $\eta \rightarrow 0$

Stochastic can be much faster when  $D$  is very large

Intermediate approach: use error over subsets of  $D$

# Maximize Conditional Log Likelihood: Gradient Ascent

$$\begin{aligned} l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) = \sum_l \ln P(Y^l | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l)) \end{aligned}$$

$$\frac{\partial l(W)}{\partial w_i} = \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

# Maximize Conditional Log Likelihood: Gradient Ascent

$$\begin{aligned} l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) = \sum_l \ln P(Y^l | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l)) \end{aligned}$$

$$\frac{\partial l(W)}{\partial w_i} = \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

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

For all  $i$ , repeat

$$w_i \leftarrow w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

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

- For MAP, need to define prior on  $W$ 
  - given  $W = \langle w_1, \dots, w_n \rangle$   
let's assume prior  $P(w_i) = N(0, \sigma)$
  - i.e., assume zero mean, Gaussian prior for each  $w_i$
- A kind of Occam's razor (simplest is best) prior
- Helps avoid very large weights and overfitting

# MAP Estimates for Logistic Regression

$$W^{MAP} = \arg \max_W P(W|Y, X) = \frac{P(Y|W, X)P(W, X)}{P(Y, X)}$$

# MAP Estimates for Logistic Regression

$$W^{MAP} = \arg \max_W P(W|Y, X) = \frac{P(Y|W, X)P(W, X)}{P(Y, X)}$$

$$= \arg \max_W P(Y|W, X)P(W, X)$$

let's assume  
 $P(W, X) = P(W)P(X)$

$$= \arg \max_W P(Y|W, X)P(W)P(X)$$

$$= \arg \max_W P(Y|W, X)P(W)$$

$$W^{MAP} = \arg \max_W [\ln P(Y|W, X) + \ln P(W)]$$

zero mean  
Gaussian  $P(W)$

$$P(W) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \sum_i w_i^2\right)$$

$$W^{MAP} = \arg \max_W [\ln P(Y|W, X) - \left(\frac{1}{2\sigma^2} \sum_i w_i^2\right)]$$

# MLE vs MAP

- Maximum conditional likelihood estimate

$$W \leftarrow \arg \max_W \ln \prod_l P(Y^l | X^l, W)$$

$$w_i \leftarrow w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

- Maximum a posteriori estimate with prior  $W \sim N(0, \sigma I)$

$$W \leftarrow \arg \max_W \ln [P(W) \prod_l P(Y^l | X^l, W)]$$

$$w_i \leftarrow w_i - \eta \lambda w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

# MAP estimates and Regularization

- Maximum a posteriori estimate with prior  $W \sim N(0, \sigma I)$

$$W \leftarrow \arg \max_W \ln [P(W) \prod_l P(Y^l | X^l, W)]$$

$$w_i \leftarrow w_i - \eta \lambda w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

called a “regularization” term

- helps reduce overfitting
- if  $P(W)$  is Gaussian, then encourages  $W$  to be near the mean of  $P(W)$  : zero here, but can easily use any mean
- used very frequently in Logistic Regression

# The Bottom Line

- 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
  - 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 ( $\pi$ )
- Then  $P(Y|X)$  is of this form, and we can directly estimate  $W$ 
$$P(Y = 1 | X = \langle X_1, \dots X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$
- Furthermore, same holds if the  $X_i$  are boolean
  - trying proving that to yourself

# What you should know:

---

- Logistic regression
  - Functional form follows from Naïve Bayes assumptions
    - For Gaussian Naïve Bayes assuming variance  $\sigma_{i,k} = \sigma_i$
    - For discrete-valued Naïve Bayes too
  - But training procedure picks parameters without making conditional independence assumption
  - MLE training: pick  $W$  to maximize  $P(Y | X, W)$
  - MAP training: pick  $W$  to maximize  $P(W | X, Y)$ 
    - ‘regularization’
    - helps reduce overfitting
- Gradient ascent/descent
  - General approach when closed-form solutions unavailable

## Extra slides

for more details, see Ng & Jordan paper

<http://www.cs.cmu.edu/~awm/10701/readings/NgJordanNIPS2001.ps>

# Generative vs. Discriminative Classifiers

Training classifiers involves estimating  $f: X \rightarrow Y$ , or  $P(Y|X)$

Generative classifiers (e.g., Naïve Bayes)

- Assume some functional form for  $P(X|Y)$ ,  $P(X)$
- Estimate parameters of  $P(X|Y)$ ,  $P(X)$  directly from training data
- Use Bayes rule to calculate  $P(Y|X=x_i)$

Discriminative classifiers (e.g., Logistic regression)

- Assume some functional form for  $P(Y|X)$
- Estimate parameters of  $P(Y|X)$  directly from training data

# Use Naïve Bayes or Logistic Regression?

Consider

- Restrictiveness of modeling assumptions
- Rate of convergence (in amount of training data) toward asymptotic hypothesis

# Naïve Bayes vs Logistic Regression

Consider  $Y$  boolean,  $X_i$  continuous,  $X = \langle X_1 \dots X_n \rangle$

Number of parameters to estimate:

- NB:

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

- LR:

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

# Naïve Bayes vs Logistic Regression

Consider  $Y$  boolean,  $X_i$  continuous,  $X = \langle X_1 \dots X_n \rangle$

Number of parameters:

- NB:  $4n + 1$
- LR:  $n+1$

Estimation method:

- NB parameter estimates are uncoupled
- LR parameter estimates are coupled

# G.Naïve Bayes vs. Logistic Regression

Recall two assumptions deriving form of LR from GNBayes:

1.  $X_i$  conditionally independent of  $X_k$  given  $Y$
2.  $P(X_i | Y = y_k) = N(\mu_{ik}, \sigma_i)$ ,  $\leftarrow$  not  $N(\mu_{ik}, \sigma_{ik})$

Consider three learning methods:

- GNB (assumption 1 only)
- GNB2 (assumption 1 and 2)
- LR

Which method works better if we have *infinite* training data, and...

- Both (1) and (2) are satisfied
- Neither (1) nor (2) is satisfied
- (1) is satisfied, but not (2)

# G.Naïve Bayes vs. Logistic Regression

[Ng & Jordan, 2002]

Recall two assumptions deriving form of LR from GNB:

1.  $X_i$  conditionally independent of  $X_k$  given  $Y$
2.  $P(X_i | Y = y_k) = N(\mu_{ik}, \sigma_i)$ ,  $\leftarrow$  not  $N(\mu_{ik}, \sigma_{ik})$

Consider three learning methods:

- GNB (assumption 1 only)
- GNB2 (assumption 1 and 2)
- LR

Which method works better if we have *infinite* training data, and...

- Both (1) and (2) are satisfied
- Neither (1) nor (2) is satisfied
- (1) is satisfied, but not (2)

# G.Naïve Bayes vs. Logistic Regression

[Ng & Jordan, 2002]

Recall two assumptions deriving form of LR from GNBayes:

1.  $X_i$  conditionally independent of  $X_k$  given  $Y$
2.  $P(X_i | Y = y_k) = N(\mu_{ik}, \sigma_i)$ ,  $\leftarrow$  not  $N(\mu_{ik}, \sigma_{ik})$

Consider three learning methods:

- GNB (assumption 1 only) -- decision surface can be non-linear
- GNB2 (assumption 1 and 2) – decision surface linear
- LR -- decision surface linear, trained without assumption 1.

Which method works better if we have *infinite* training data, and...

- Both (1) and (2) are satisfied:  $LR = GNB2 = GNB$
- (1) is satisfied, but not (2) :  $GNB > GNB2, GNB > LR, LR > GNB2$
- Neither (1) nor (2) is satisfied:  $GNB > GNB2, LR > GNB2, LR > <GNB$

# G.Naïve Bayes vs. Logistic Regression

[Ng & Jordan, 2002]

What if we have only finite training data?

They converge at different rates to their asymptotic ( $\infty$  data) error

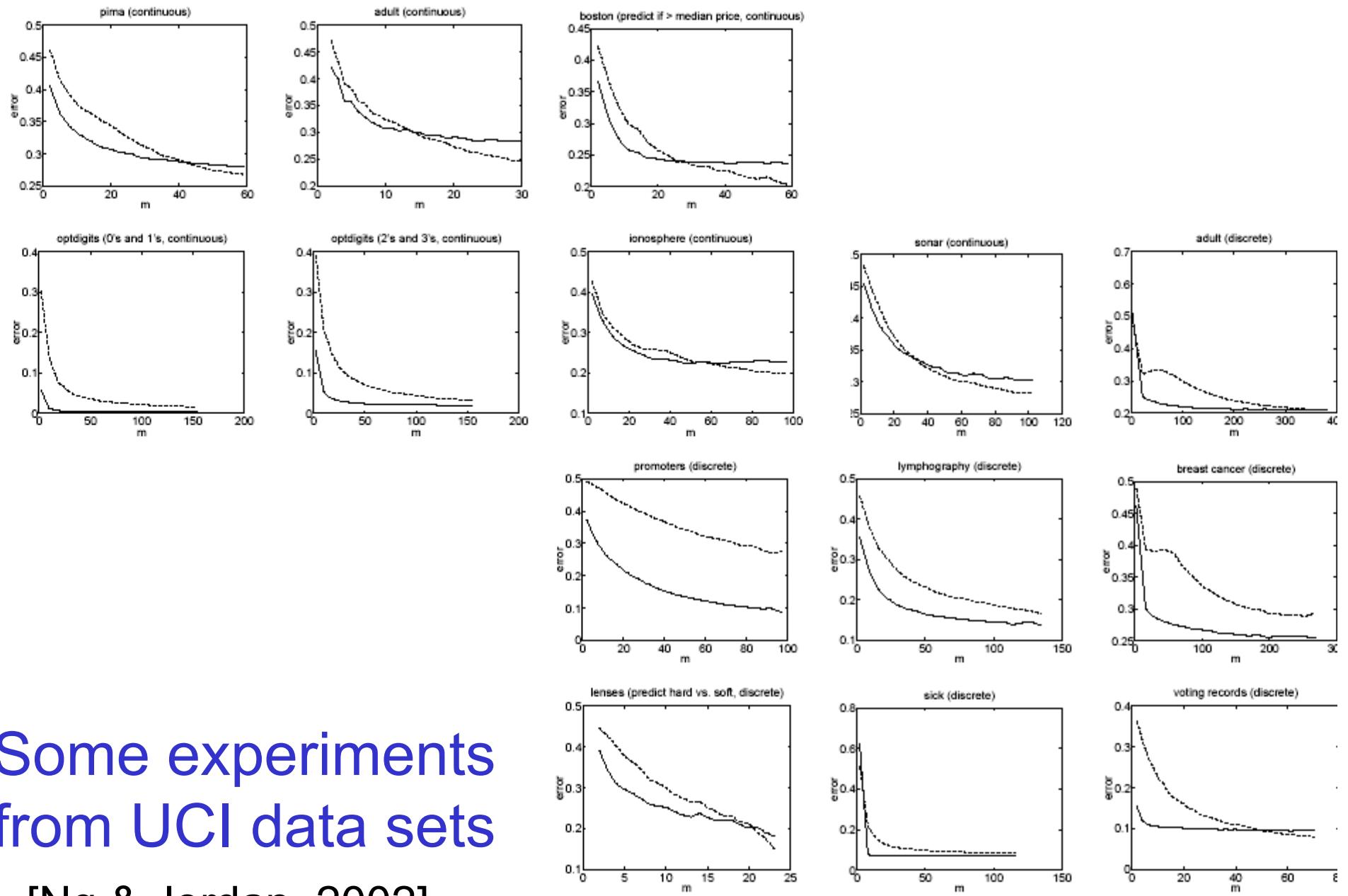
Let  $\epsilon_{A,n}$  refer to expected error of learning algorithm A after  $n$  training examples

Let  $d$  be the number of features:  $\langle X_1 \dots X_d \rangle$

$$\epsilon_{LR,n} \leq \epsilon_{LR,\infty} + O\left(\sqrt{\frac{d}{n}}\right)$$

$$\epsilon_{GNB,n} \leq \epsilon_{GNB,\infty} + O\left(\sqrt{\frac{\log d}{n}}\right)$$

So, GNB requires  $n = O(\log d)$  to converge, but LR requires  $n = O(d)$



# Some experiments from UCI data sets

[Ng & Jordan, 2002]

Figure 1: Results of 15 experiments on datasets from the UCI Machine Learning repository. Plots are of generalization error vs.  $m$  (averaged over 1000 random train/test splits). Dashed line is logistic regression; solid line is naive Bayes.

# Naïve Bayes vs. Logistic Regression

The bottom line:

GNB2 and LR both use linear decision surfaces, GNB need not

Given infinite data, LR is better or equal to GNB2 because *training procedure* does not make assumptions 1 or 2 (though our derivation of the form of  $P(Y|X)$  did).

But GNB2 converges more quickly to its perhaps-less-accurate asymptotic error

And GNB is both more biased (assumption1) and less (no assumption 2) than LR, so either might outperform the other

# What you should know:

---

- Logistic regression
  - Functional form follows from Naïve Bayes assumptions
    - For Gaussian Naïve Bayes assuming variance  $\sigma_{i,k} = \sigma_i$
    - For discrete-valued Naïve Bayes too
  - But training procedure picks parameters without making conditional independence assumption
  - MLE training: pick  $W$  to maximize  $P(Y | X, W)$
  - MAP training: pick  $W$  to maximize  $P(W | X, Y)$ 
    - ‘regularization’
    - helps reduce overfitting
- Gradient ascent/descent
  - General approach when closed-form solutions unavailable
- Generative vs. Discriminative classifiers
  - Bias vs. variance tradeoff