# CSE546: Logistic Regression
## Winter 2012

Luke Zettlemoyer

Slides adapted from Carlos Guestrin

# Lets take a(nother) probabilistic approach!!!

- **Previously: directly estimate the data distribution P(X,Y)!**
  - challenging due to size of distribution!
  - make Naïve Bayes assumption: only need $P(X_i|Y)$!
- **But wait, we classify according to:**
  - $\max_Y P(Y|X)$
- **Why not learn P(Y|X) directly?**

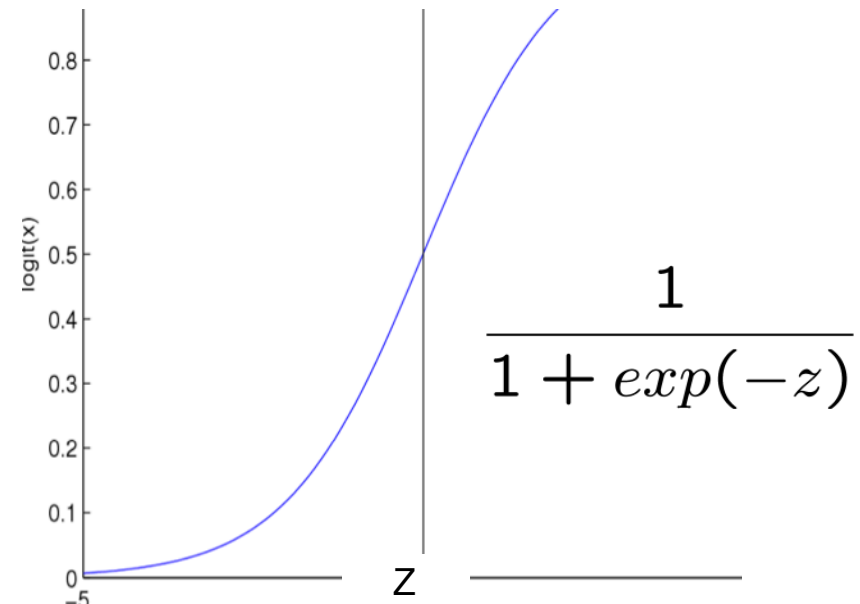| mpg | cylinders | displacemen | horsepower | weight | acceleration | modelyear | make |
|-----|-----------|-------------|------------|--------|--------------|-----------|------|
| | | | | | | | |
| good | 4 | 97 | 75 | 2265 | 18.2 | 77 | asia |
| bad | 6 | 199 | 90 | 2648 | 15 | 70 | ameri |
| bad | 4 | 121 | 110 | 2600 | 12.8 | 77 | europ |
| bad | 8 | 350 | 175 | 4100 | 13 | 73 | ameri |
| bad | 6 | 198 | 95 | 3102 | 16.5 | 74 | ameri |
| bad | 4 | 108 | 94 | 2379 | 16.5 | 73 | asia |
| bad | 4 | 113 | 95 | 2228 | 14 | 71 | asia |
| bad | 8 | 302 | 139 | 3570 | 12.8 | 78 | ameri |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| good | 4 | 120 | 79 | 2625 | 18.6 | 82 | ameri |
| bad | 8 | 455 | 225 | 4425 | 10 | 70 | ameri |
| good | 4 | 107 | 86 | 2464 | 15.5 | 76 | europ |
| bad | 5 | 131 | 103 | 2830 | 15.9 | 78 | europ |
| | | | | | | | |

# Logistic Regression

- ## Learn P(Y|**X**) directly!

  - Assume a particular functional form

  - Sigmoid applied to a linear function of the data:

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

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

$$P(Y = 0|X) = \frac{\exp(w_0 + \sum_{i=1}^{n} w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^{n} w_i X_i)}$$

**Features can be discrete or continuous!**

# Logistic Regression: decision boundary

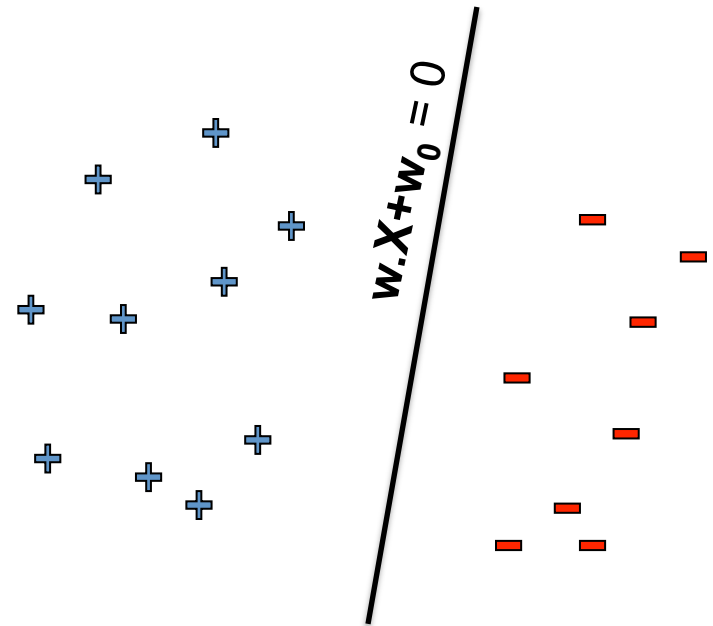$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^{n} w_i X_i)} \qquad P(Y = 0|X) = \frac{\exp(w_0 + \sum_{i=1}^{n} w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^{n} w_i X_i)}$$

- Prediction: Output the Y with highest P(Y|X)
  - For binary Y, output Y=0 if

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

$$1 < \exp(w_0 + \sum_{i=1}^{n} w_i X_i)$$

$$0 < w_0 + \sum_{i=1}^{n} w_i X_i$$

w.X+w$_0$ = 0

A Linear Classifier!

# Logistic regression for discrete classification

Logistic regression in more general case, where set of possible *Y* is {$y_1,...,y_R$}

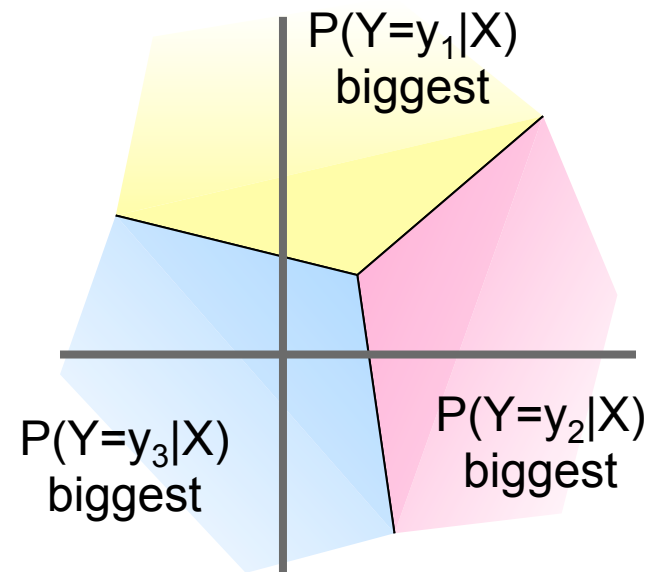- Define a weight vector $w_i$ for each $y_i$, i=1,…,R-1

$$P(Y = 1|X) \propto \exp\left(w_{10} + \sum_i w_{1i}X_i\right)$$

$$P(Y = 2|X) \propto \exp\left(w_{20} + \sum_i w_{2i}X_i\right)$$

...

$$P(Y = r|X) = 1 - \sum_{j=1}^{r-1} P(Y = j|X)$$

P(Y=$y_1$|X) biggest

P(Y=$y_2$|X) biggest

P(Y=$y_3$|X) biggest

# Logistic regression: discrete Y

- Logistic regression in more general case, where
  $Y$ is in the set $\{y_1,...,y_R\}$

  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!**

# Loss functions / Learning Objectives: Likelihood v. Conditional Likelihood

- Generative (Naïve Bayes) Loss function:
  **Data likelihood**

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

- But, discriminative (logistic regression) loss function:
  **Conditional Data Likelihood**

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

  – Doesn't waste effort learning P(X) – focuses on P(Y|**X**) all that matters for classification
  – Discriminative models cannot compute P($\mathbf{x}^j$|**w**)!

# Conditional Log Likelihood
## (the binary case only)

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

$$l(\mathbf{w}) \equiv \sum_j \ln P(y^j|\mathbf{x}^j, \mathbf{w}) \qquad P(Y = 1|\mathbf{X}, \mathbf{w}) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

equal because $y^j$ is in {0,1}

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

remaining steps: substitute definitions, expand logs, and simplify

$$= \sum_j y^j \ln \frac{e^{w_0 + \sum_i w_i X_i}}{1 + e^{w_0 + \sum_i w_i X_i}} + (1 - y^j) \ln \frac{1}{1 + e^{w_0 + \sum_i w_i X_i}}$$

$$\cdots$$

# Logistic Regression Parameter Estimation: Maximize Conditional Log Likelihood

$$
\begin{aligned}
l(\mathbf{w}) &\equiv \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) \\
&= \sum_j y^j \left( w_0 + \sum_i^n w_i x_i^j \right) - \ln\left( 1 + exp\left( w_0 + \sum_i^n w_i x_i^j \right) \right)
\end{aligned}
$$

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

# Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave !

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

**Learning rate, $\eta > 0$**

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

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

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

# Maximize Conditional Log Likelihood: Gradient ascent

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

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

$$\frac{\partial l(w)}{\partial w_i} = \sum_j \left[ \frac{\partial}{\partial w} y^j (w_0 + \sum_i w_i x_i^j) - \frac{\partial}{\partial w} \ln \left( 1 + \exp(w_0 + \sum_i w_i x_i^j) \right) \right]$$

$$= \sum_j \left[ y^j x_i^j - \frac{x_i^j \exp(w_0 + \sum_i w_i x_i^j)}{1 + \exp(w_0 + \sum_i w_i x_i^j)} \right]$$

$$= \sum_j x_i^j \left[ y^j - \frac{\exp(w_0 + \sum_i w_i x_i^j)}{1 + \exp(w_0 + \sum_i w_i x_i^j)} \right]$$

$$\frac{\partial l(w)}{\partial w_i} = \sum_j x_i^j \left( y^j - P(Y^j = 1 | x^j, w) \right)$$

# Gradient Descent for LR

Gradient ascent algorithm: (learning rate $\eta > 0$)

`do:`

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta \sum_j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w})]$$

`For i=1…n: (iterate over weights)`

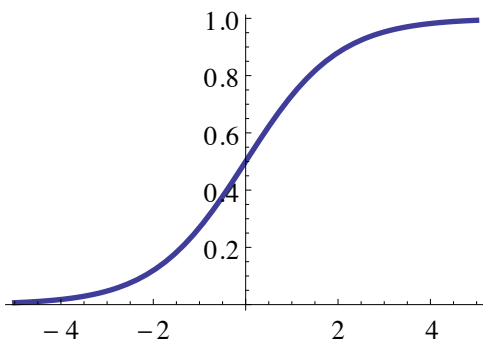$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w})]$$

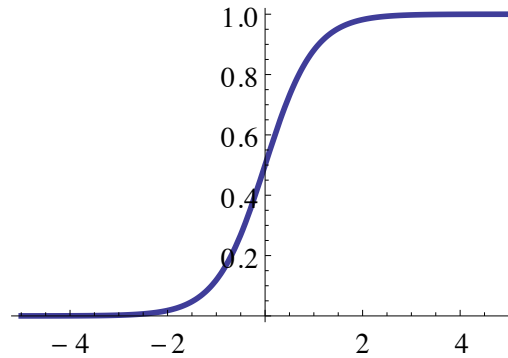`until "change" < ε`

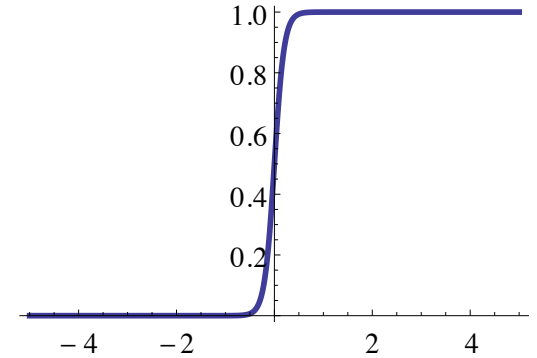Loop over training examples!

# Large parameters…

$$\frac{1}{1 + e^{-ax}}$$



a=1    a=5    a=10

- **Maximum likelihood solution: prefers higher weights**
  - higher likelihood of (properly classified) examples close to decision boundary
  - larger influence of corresponding features on decision
  - *can cause overfitting!!!*
- **Regularization: penalize high weights**
  - again, more on this later in the quarter

# 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 **w**
  - Normal distribution, zero mean, identity covariance
  - "Pushes" parameters towards zero

$$p(\mathbf{w}) = \prod_i \frac{1}{\kappa\sqrt{2\pi}} \; e^{\frac{-w_i^2}{2\kappa^2}}$$

- Often called ***Regularization***
  - Helps avoid very large weights and overfitting

- MAP estimate:

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

# M(C)AP as Regularization

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

- Add log p(w) to objective:

$$\ln p(w) \propto -\frac{\lambda}{2} \sum_i w_i^2 \qquad \frac{\partial \ln p(w)}{\partial w_i} = -\lambda w_i$$

  – Quadratic penalty: drives weights towards zero
  – Adds a negative linear term to the gradients

**Penalizes high weights, also applicable in linear regression**

# MLE vs. MAP

- Maximum conditional likelihood estimate

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

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

- Maximum conditional a posteriori estimate

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

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

# Logistic regression v. Naïve Bayes

- Consider learning f: X → Y, where
  - X is a vector of real-valued features, < $X_1$ ... $X_n$ >
  - Y is boolean
- Could use a Gaussian Naïve Bayes classifier
  - assume all $X_i$ are conditionally independent given Y
  - model $P(X_i \mid Y = y_k)$ 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)?

$$P(Y = 1 | X =< X_1, ...X_n >) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

**Cool!!!!**

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

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

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

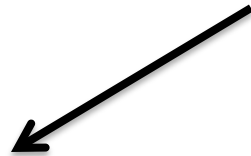$$= \frac{1}{1 + \exp(\ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)})}$$

up to now, all arithmetic

only for Naïve Bayes models

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

Looks like a setting for $w_0$?

Can we solve for $w_i$ ?
- Yes, but only in Gaussian case

# Ratio of class-conditional probabilities

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

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

$$= \ln \left[ \frac{\frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x_i - \mu_{i0})^2}{2\sigma_i^2}}}{\frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2}}} \right]$$
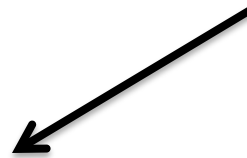
$$= -\frac{(x_i - \mu_{i0})^2}{2\sigma_i^2} + \frac{(x_i - \mu_{i1})^2}{2\sigma_i^2}$$

$$\dots$$

$$= \frac{\mu_{i0} + \mu_{i1}}{\sigma_i^2} x_i + \frac{\mu_{i0}^2 + \mu_{i1}^2}{2\sigma_i^2}$$

Linear function!
Coefficents
expressed with
original Gaussian
parameters!

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

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

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

$$\boxed{\sum_i \left( \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right)}$$
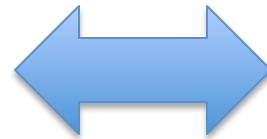
$$P(Y = 1 | X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

$$w_0 = \ln \frac{1 - \theta}{\theta} + \frac{\mu_{i0}^2 + \mu_{i1}^2}{2\sigma_i^2}$$

$$w_i = \frac{\mu_{i0} + \mu_{i1}}{\sigma_i^2}$$

# Gaussian Naïve Bayes vs. Logistic Regression

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

**Can go both ways, we only did one way**

**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** $P(\mathbf{X}|Y)$ **in learning!!!**
- **Loss function!!!**
  - Optimize different functions ! Obtain different solutions

# Naïve Bayes vs. Logistic Regression

Consider Y boolean, $X_i$ continuous, $X=<X_1 \ldots X_n>$

## Number of parameters:

- Naïve Bayes: $4n + 1$
- Logistic Regression: $n+1$

## Estimation method:

- Naïve Bayes parameter estimates are uncoupled
- Logistic Regression parameter estimates are coupled

# Naïve Bayes vs. Logistic Regression

[Ng & Jordan, 2002]

- Generative vs. Discriminative classifiers

- Asymptotic comparison

  (# training examples → infinity)

  – when model correct

    - GNB (with class independent variances) and LR produce identical classifiers

  – when model incorrect

    - LR is less biased – does not assume conditional independence

      – **therefore LR expected to outperform GNB**

# Naïve Bayes vs. Logistic Regression

[Ng & Jordan, 2002]

- Generative vs. Discriminative classifiers

- Non-asymptotic analysis
  - convergence rate of parameter estimates,
    - (n = # of attributes in X)
    - Size of training data to get close to infinite data solution
    - Naïve Bayes needs $O(\log n)$ samples
    - Logistic Regression needs $O(n)$ samples

  - GNB converges more quickly to its (perhaps less helpful) asymptotic estimates

Naïve bayes
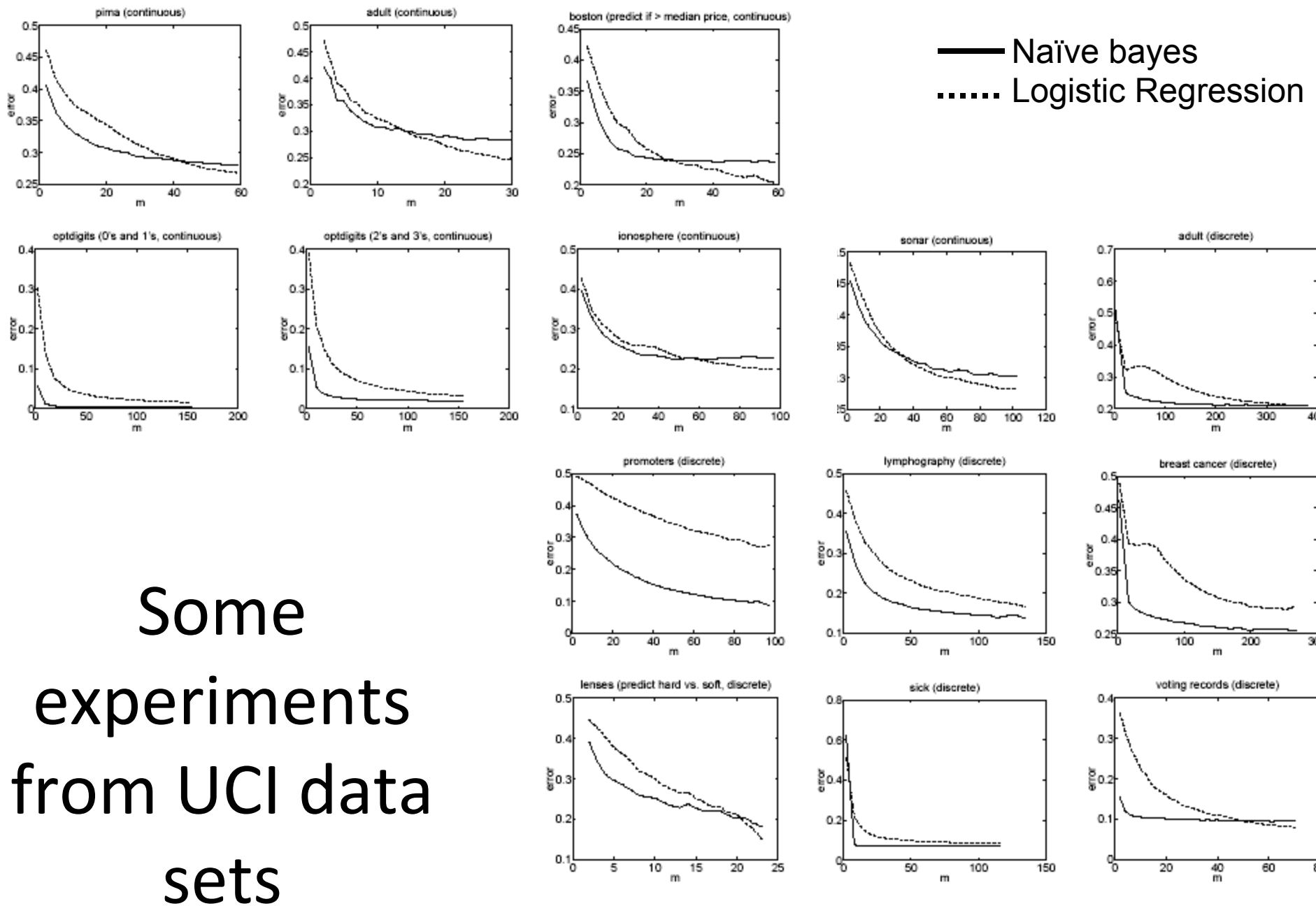...... Logistic Regression

Some experiments from UCI data sets

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.

# 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 ! assumption on P(**X**|Y)
  - LR: Functional form of P(Y|**X**), no assumption on P(**X**|Y)
- LR is a linear classifier
  - decision rule is a hyperplane
- LR optimized by conditional likelihood
  - no closed-form solution
  - concave ! global optimum with gradient ascent
  - Maximum conditional a posteriori corresponds to regularization
- Convergence rates
  - GNB (usually) needs less data
  - LR (usually) gets to better solutions in the limit