

# Logistic Regression

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

September 28<sup>th</sup>, 2009

©Carlos Guestrin 2005-2009

1

## Logistic Regression

Logistic function (or Sigmoid):  $\frac{1}{1 + \exp(-z)}$

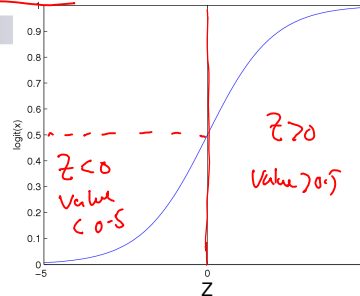
■ Learn  $P(Y|X)$  directly!

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

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

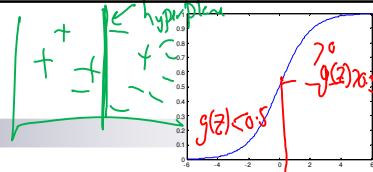
regression  $\rightarrow \mathbb{R}$   
 Logistic regression  $\rightarrow \{0,1\}$   
 Params:  $w_0, w_1, \dots, w_n$   
 like applying logistic regression to linear regression function

**Features can be discrete or continuous!**



2

## Logistic Regression – a Linear classifier



$X = \{x_1, \dots, x_n\}$  (each  $x_i$  can be anything)  
 $P(Y=1|X, w) < 0.5$   
 binary class  
 predict  $Y=0$   
 $w_0 + \sum_i w_i x_i > 0$   
 $w_0 + \sum_i w_i x_i = 0$  ← hyper plane  
 $w_0 + \sum_i w_i x_i < 0$   
 $P(Y=1|X, w) > 0.5$   
 predict  $Y=1$

©Carlos Guestrin 2005-2009

3

## Logistic regression for more than 2 classes

- Logistic regression in more general case, where

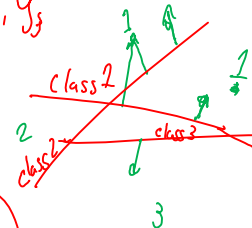
$Y \in \{y_1, \dots, y_R\}$   $R$  classes:  $y_1, y_2, y_3$

$$P(Y=1|X) \propto e^{w_{10} + \sum_i w_{1i} x_i}$$

$$P(Y=2|X) \propto e^{w_{20} + \sum_i w_{2i} x_i}$$

$\vdots$

$$P(Y=3|X) = 1 - P(Y=1|X) - P(Y=2|X)$$



©Carlos Guestrin 2005-2009

4

# Logistic regression more generally

- Logistic regression in more general case, where  $Y \in \{y_1, \dots, y_R\}$

for  $k < R$

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

*learn  $w_{ki}$*

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!**

©Carlos Guestrin 2005-2009

5

# Loss functions: Likelihood v. Conditional Likelihood

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

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

*vector of features class*  
*conditional likelihood*

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

**Conditional Data Likelihood**

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

*only this part*

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

©Carlos Guestrin 2005-2009

6

## Expressing Conditional Log Likelihood

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

$$= \sum_j \begin{cases} \ln P(Y=0 | \mathbf{x}^j, \mathbf{w}) & \text{if } y^j=0 \\ \ln P(Y=1 | \mathbf{x}^j, \mathbf{w}) & \text{if } y^j=1 \end{cases}$$

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

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

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

$$= \sum_j y_j \left[ \ln \frac{e^{w_0 + \sum_i w_i x_i}}{1 + e^{w_0 + \sum_i w_i x_i}} \right] + (1 - y_j) \left[ \ln \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}} \right]$$

$$= \sum_j y_j (w_0 + \sum_i w_i x_i) - y_j \ln(1 + e^{w_0 + \sum_i w_i x_i}) - (1 - y_j) \ln(1 + e^{w_0 + \sum_i w_i x_i})$$

©Carlos Guestrin 2005-2009

7

## Maximizing Conditional Log Likelihood

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

$$= \sum_{j=1}^N y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))$$

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

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

Multiclass  $\rightarrow$  still concave, but eqn is slightly longer...

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

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

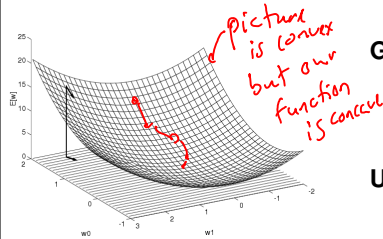
**Good news:** concave functions 'easy' to optimize

©Carlos Guestrin 2005-2009

8

# Optimizing concave function – Gradient ascent

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



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

**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)

©Carlos Guestrin 2005-2009

9

## Maximize Conditional Log Likelihood: Gradient ascent

$$\frac{\partial}{\partial x} \ln f(x) = \frac{f'(x)}{f(x)}$$

$$\frac{\partial}{\partial x} e^{ax} = a e^{ax}$$

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

$$\begin{aligned} \frac{\partial l(\mathbf{w})}{\partial w_i} &= \sum_j \underbrace{\frac{\partial}{\partial w_i} y^j (w_0 + \sum_i w_i x_i^j)}_{y^j x_i^j} - \underbrace{\frac{\partial}{\partial w_i} \ln(1 + e^{w_0 + \sum_i w_i x_i^j})}_{\frac{x_i^j e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}}} \end{aligned}$$

$$= \sum_j x_i^j \left( y^j - \underbrace{\frac{e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}}}_{P(Y=1 | \mathbf{x}^j, \mathbf{w})} \right)$$

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

©Carlos Guestrin 2005-2009

10

# Gradient Descent for LR

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

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

For  $i=1, \dots, n$ ,

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

repeat

©Carlos Guestrin 2005-2009

11

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

$$p(\mathbf{w} | Y, \mathbf{X}) \propto \underbrace{P(Y | \mathbf{X}, \mathbf{w})}_{\text{likelihood}} \underbrace{p(\mathbf{w})}_{\text{prior over params}}$$

- 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
- More on this later in the semester

- MAP estimate

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

©Carlos Guestrin 2005-2009

12

# M(C)AP as Regularization

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

$= \ln p(\mathbf{w}) + \sum_j \ln P(y^j | \mathbf{x}^j, \mathbf{w})$

$= \ln \prod_i \frac{1}{\kappa \sqrt{2\pi}} e^{-\frac{w_i^2}{2\kappa^2}}$

$= \sum_i \ln \frac{1}{\kappa \sqrt{2\pi}} + \sum_i -\frac{w_i^2}{2\kappa^2}$

const. doesn't depend on  $\mathbf{w} \rightarrow$  ignore

$\Rightarrow \sum_j \ln P(y^j | \mathbf{x}^j, \mathbf{w}) - \sum_i \frac{w_i^2}{2\kappa^2}$

as before penalty for large  $w_i$

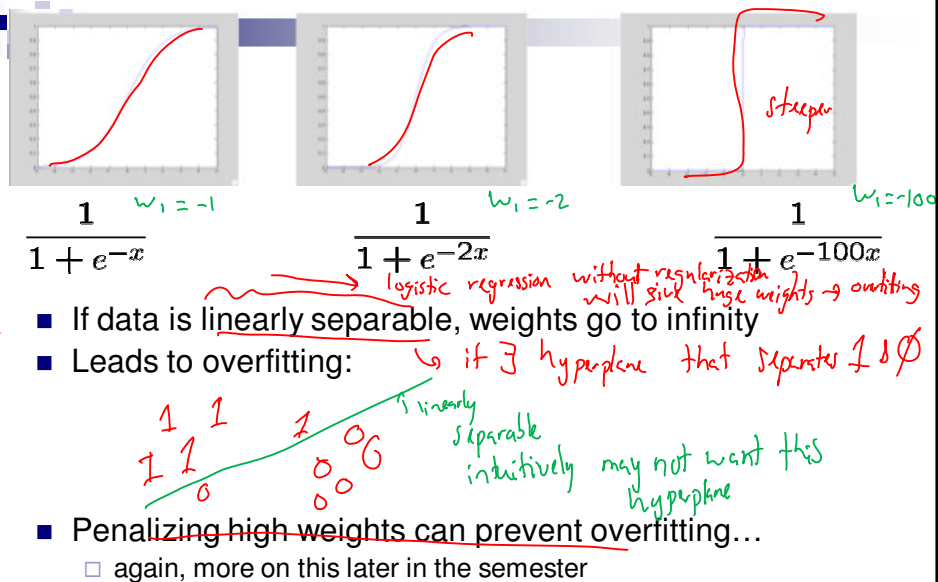
$p(\mathbf{w}) = \prod_i \frac{1}{\kappa \sqrt{2\pi}} e^{-\frac{w_i^2}{2\kappa^2}}$

Zero mean gaussian prior

variance  $\kappa^2$

Penalizes high weights, also applicable in linear regression

## Large parameters $\rightarrow$ Overfitting



## Gradient of M(C)AP

$$\ln ab = \ln a + \ln b$$

$$\frac{\partial}{\partial w_i} w_i^2 = 2w_i$$

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

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

$$= \underbrace{\frac{\partial}{\partial w_i} \ln p(\mathbf{w})}_{\text{extra term}} + \underbrace{\frac{\partial}{\partial w_i} \ln \left( \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right)}_{\text{as before}}$$

$$\frac{\partial}{\partial w_i} \left[ \text{const.} + \sum_i -\frac{w_i^2}{2\kappa^2} \right] = -\frac{w_i}{\kappa^2} \leftarrow \text{extra term pushes } w_i \text{ towards } 0$$

if  $w_i > 0 \Rightarrow -\frac{w_i}{\kappa^2}$  is very negative } push  $w_i$  to 0  
 $w_i < 0 \Rightarrow -\frac{w_i}{\kappa^2}$  is very positive } push  $w_i$  to 0

©Carlos Guestrin 2005-2009

15

## MLE vs MAP

in practice, ALWAYS REGULARIZE YOUR LR

- Maximum conditional likelihood estimate

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

Typically don't regularize  $w_0$

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

update rule gradient ascent (regularized)

- Maximum conditional a posteriori estimate

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

extra term  $\lambda = \frac{1}{\kappa^2}$

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

©Carlos Guestrin 2005-2009

16



# 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!!!!

©Carlos Guestrin 2005-2009

17

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

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

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

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

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

reminds  
me of  
 $w_0$

reminds me of  
 $w_i X_i$

only assumption  
thus far is  
NB

let's check

©Carlos Guestrin 2005-2009

18

## Ratio of class-conditional probabilities

$$\ln \frac{a}{b} = \ln a - \ln b$$

$$\ln e^{\text{Something}} = \text{Something}$$

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

$$\ln \frac{\frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x_i - \mu_{i0})^2}{2\sigma_i^2}}}{\frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2}}}$$

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

$$= -\cancel{x_i^2} + 2\mu_{i0}x_i - \mu_{i0}^2 + \cancel{x_i^2} - 2\mu_{i1}x_i + \mu_{i1}^2$$

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

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

formula for GNB,

$\sigma_i$  class indep

goal.  
 $\ln \frac{P(X_i | Y=0)}{P(X_i | Y=1)}$   
 to look like  
 $w_i x_i$

part of  $w_0$

©Carlos Guestrin 2005-2009

19

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

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

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

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

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

$$w_0 = \ln \left( \frac{1-\theta}{\theta} \right) + \sum_i \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}$$

©Carlos Guestrin 2005-2009

20

## Gaussian Naïve Bayes v. Logistic Regression

- 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  $\Rightarrow$  Obtain different solutions

$LR: P(Y|X, w)$   
 $NB: P(Y, X|w)$

©Carlos Guestrin 2005-2009

21

## 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

©Carlos Guestrin 2005-2009

22

## G. Naïve Bayes vs. Logistic Regression 1

[Ng & Jordan, 2002]

### ■ Generative and Discriminative classifiers

### ■ Asymptotic comparison (# training examples $\rightarrow$ infinity)

- when <sup>GNB independent</sup> model correct
  - GNB (with class independent variances) and LR produce identical classifiers
- when <sup>GNB</sup> model incorrect
  - LR is less biased – does not assume conditional independence
    - **therefore LR expected to outperform GNB**

©Carlos Guestrin 2005-2009

23

## G. Naïve Bayes vs. Logistic Regression 2

[Ng & Jordan, 2002]

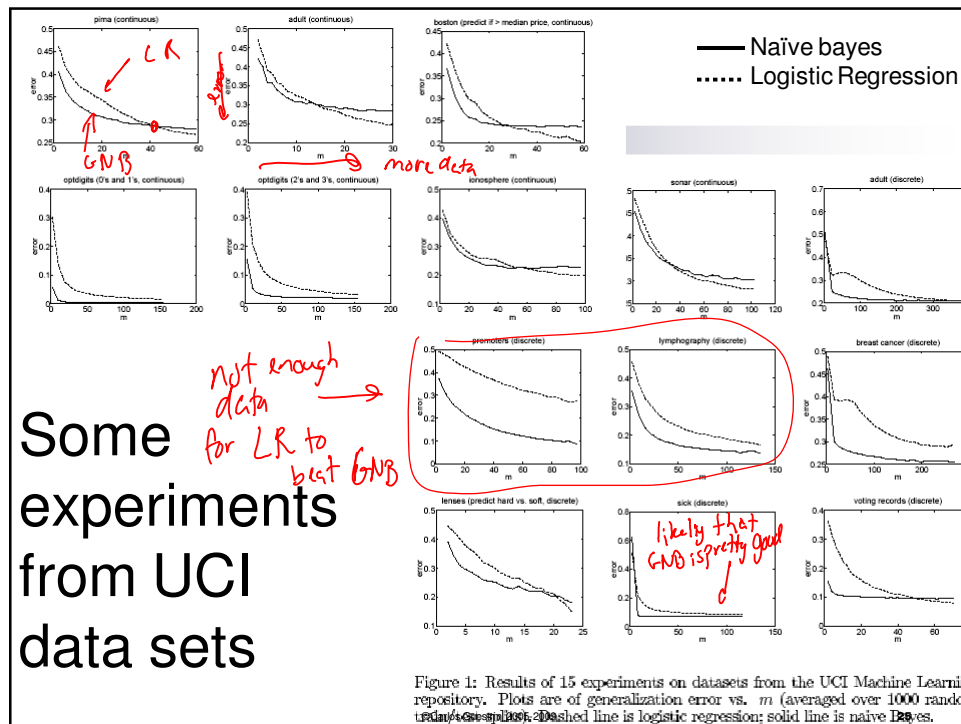
### ■ Generative and Discriminative classifiers

### ■ Non-asymptotic analysis

- convergence rate of parameter estimates,  $n = \#$  of attributes in  $X$ 
  - Size of training data to get close to infinite data solution
  - GNB needs  $O(\log n)$  samples *n features (because of indep.)*
  - LR needs  $O(n)$  samples *need more data*
- **GNB converges more quickly to its (perhaps less helpful) asymptotic estimates**

©Carlos Guestrin 2005-2009

24



## 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  $P(\mathbf{X}|\mathbf{Y})$
  - LR: Functional form of  $P(\mathbf{Y}|\mathbf{X})$ , no assumption on  $P(\mathbf{X}|\mathbf{Y})$
- 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