

# Learning Logistic Regressors by Gradient Descent

Machine Learning – CSE446

Carlos Guestrin

University of Washington

April 17, 2013

©Carlos Guestrin 2005-2013

1

## Classification

### ■ Learn: $h: \mathbf{X} \mapsto Y$

- $\mathbf{X}$  – features
- $Y$  – target classes

$\mathbf{X} \equiv (\text{GPA}, \text{ML grade})$   
in reg.:  $y = \text{salary}$  (cont. value)

now:  $y$  is discrete

e.g.  $Y \equiv \{\text{hired}, \text{not hired}\}$

### ■ Conditional probability: $P(Y|\mathbf{X})$

$P(y = \text{hired} \mid \text{GPA} = 3.6, \text{ML grade} = 3.9)$

### ■ Suppose you know $P(Y|\mathbf{X})$ exactly, how should you classify?

- Bayes optimal classifier:

$$\hat{y} = \arg \max_y P(Y=y \mid X=x)$$

### ■ How do we estimate $P(Y|\mathbf{X})$ ?

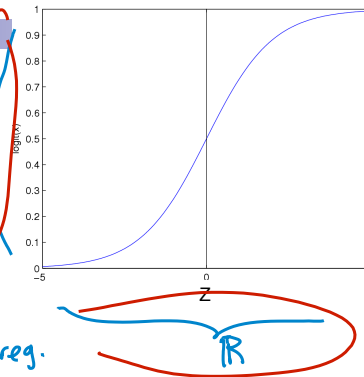
$P(\text{hired} \mid 3.6, 3.9) = 0.8$   
 $P(\text{not hired} \mid 3.6, 3.9) = 0.2$   
 $\Rightarrow$  predict hired

©Carlos Guestrin 2005-2013

2

# Logistic Regression

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



Learn  $P(Y|X)$  directly

- Assume a particular functional form for link function
- Sigmoid applied to a linear function of the input features:

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

choice

$z$ : linear just like in reg.

$w_0 + \sum_i w_i x_i$  ← not bounded, could be neg.

after logistic fcn, output is in  $[0,1]$

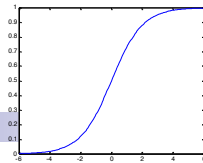
**Features can be discrete or continuous!**

©Carlos Guestrin 2005-2013

3

## Logistic Regression – a Linear classifier

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



$$P(Y=0|X, w) =$$

$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

$w_0 + \sum_i w_i x_i > 0$   
 $\Rightarrow g(w_0 + \sum_i w_i x_i) < 0.5$   
 $\Rightarrow P(Y=0|X, w) < 0.5$   
 $\Rightarrow \text{predict class} = 1$

$$w_0 + \sum_i w_i x_i = 0$$

$w_0 + \sum_i w_i x_i < 0$   
 $\Rightarrow g(w_0 + \sum_i w_i x_i) > 0.5$   
 $\Rightarrow \text{predict class} = 0$

©Carlos Guestrin 2005-2013

4

# Loss function: Conditional Likelihood

- Have a bunch of iid data of the form:

$(x^j, y^j)_{j=1}^N = D = (D_X, D_Y)$   
iid

$x$  (6PA:3.2)  
 $y$  (hired)

- Discriminative (logistic regression) loss function:

Conditional Data Likelihood

$\arg \max_w P(D_Y | D_X, w) = \arg \max_x \prod_{j=1}^N P(y^j | x^j, w)$   
 $= \arg \max_w \ln \prod_{j=1}^N P(y^j | x^j, w) = \arg \max_x \sum_{j=1}^N \ln P(y^j | x^j, w)$

$$\ln P(D_Y | D_X, w) = \sum_{j=1}^N \ln P(y^j | x^j, w)$$

©Carlos Guestrin 2005-2013

5

## Maximizing Conditional Log Likelihood

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

$$= \sum_j y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))$$

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

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

Good news:  $l(w)$  is concave function of  $w$ , no local optima problems

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

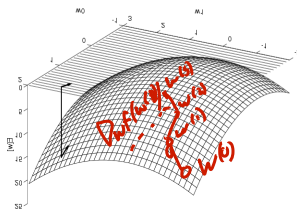
Good news: concave functions easy to optimize

©Carlos Guestrin 2005-2013

6

# Optimizing concave function – Gradient ascent

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



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

Step size,  $\eta > 0$

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

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

For Hw 2  
 $\eta$  will be constant

- Gradient ascent is simplest of optimization approaches
  - e.g., Conjugate gradient ascent can be much better

Often, especially in proof,  $\eta$  gets smaller with iterations  
e.g.  $\eta_t = \frac{\alpha}{t}$   $\alpha \leftarrow \text{constant}$

©Carlos Guestrin 2005-2013

7

## Maximize Conditional Log Likelihood: Gradient ascent

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

$\frac{\partial e^{f(s)}}{\partial s} = f'(s) e^{f(s)}$

$$l(\mathbf{w}) = \sum_{j=1}^N y^j (w_0 + \sum_{i=1}^k w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_{i=1}^k w_i x_i^j))$$

Der:  $\frac{\partial l}{\partial w_i} = \sum_{j=1}^N \left[ y^j x_i^j - \frac{x_i^j \exp(w_0 + \sum_{i=1}^k w_i x_i^j)}{1 + \exp(w_0 + \sum_{i=1}^k w_i x_i^j)} \right]$

$\hat{p}(y=1 | x_i^j, \mathbf{w})$

$$\frac{\partial l}{\partial w_i} = \sum_{j=1}^N x_i^j (y^j - \hat{p}(y=1 | x_i^j, \mathbf{w}))$$

weighted by contribution of  $i$ th feature to point  $j$

how far is my prediction from the truth

©Carlos Guestrin 2005-2013

8

# Gradient Ascent for LR

*Start from some  $w^{(0)}$  e.g.  $\emptyset$*

*revisit soon*

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

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

*step size*

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

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

repeat

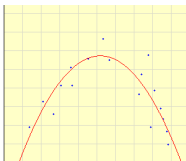
©Carlos Guestrin 2005-2013

9

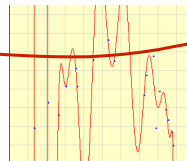
# Regularization in linear regression

- Overfitting usually leads to very large parameter choices, e.g.:

$$-2.2 + 3.1 X - 0.30 X^2$$



$$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$$

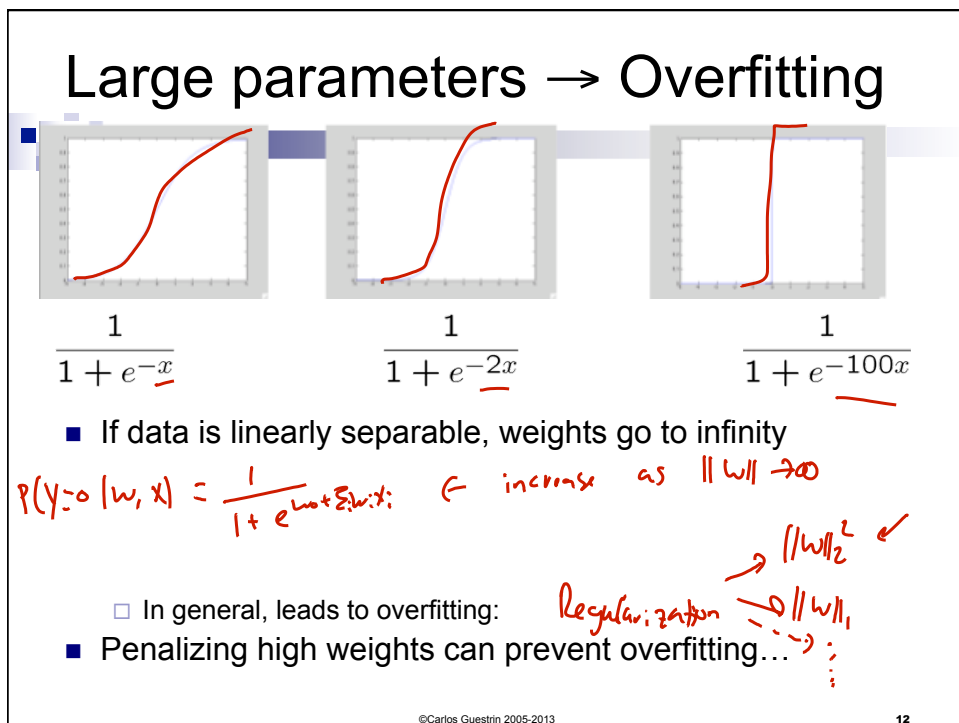
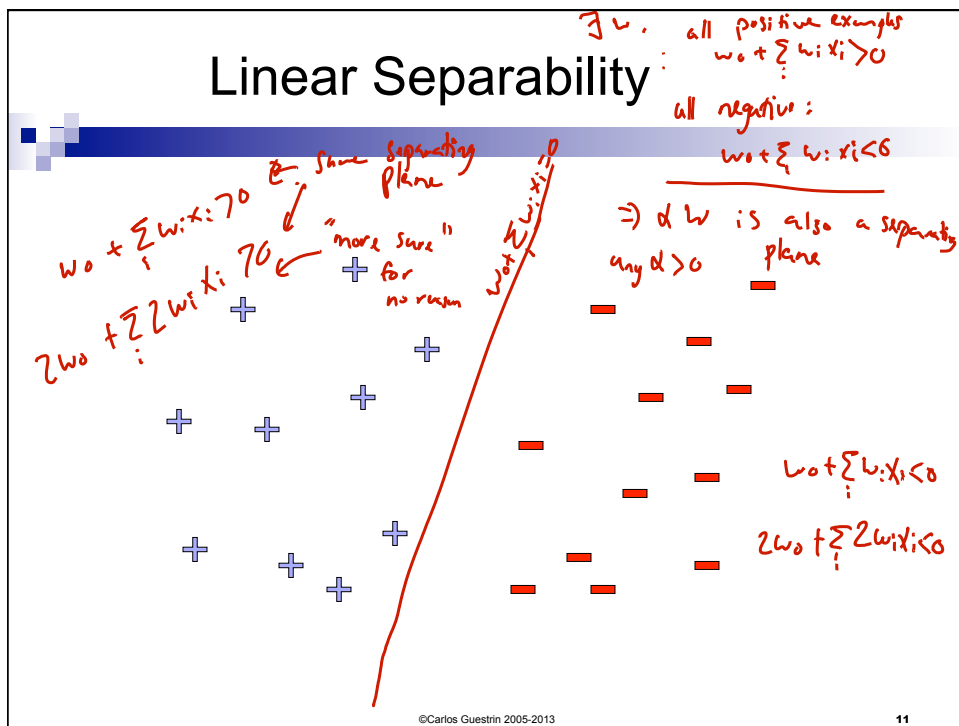


- Regularized least-squares (a.k.a. ridge regression), for  $\lambda > 0$ :

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 + \lambda \sum_{i=1}^k w_i^2$$

©Carlos Guestrin 2005-2013

10



## Regularized Conditional Log Likelihood

- Add regularization penalty, e.g.,  $L_2$ :

$$\ell(\mathbf{w}) = \ln \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

*Handwritten notes:*  $\sum_{i=1}^k w_i^2$  (under the regularization term)

- Practical note about  $w_0$ :

*Handwritten note:* don't regularize

- Gradient of regularized likelihood:

$$\frac{\partial \ell}{\partial w_i} = \frac{\partial}{\partial w_i} \left[ \ln \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right] - \frac{\lambda}{2} \frac{\partial \|\mathbf{w}\|_2^2}{\partial w_i}$$

*Handwritten notes:*  $\lambda w_i$  (above the second term),  $\frac{\partial \|\mathbf{w}\|_2^2}{\partial w_i}$  (under the second term)

©Carlos Guestrin 2005-2013

13

## Standard v. Regularized Updates

- Maximum conditional likelihood estimate

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

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

*Handwritten note:* without regularization

- Regularized maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) - \frac{\lambda}{2} \sum_{i=1}^k w_i^2$$

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

*Handwritten note:* push towards 0 (above the  $-\lambda w_i^{(t)}$  term)

©Carlos Guestrin 2005-2013

14

*w\* is optimal solution to learning problem*

## Please Stop!! Stopping criterion

$$\ell(\mathbf{w}) = \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) - \lambda \|\mathbf{w}\|_2^2$$

- When do we stop doing gradient descent? *ascend or user-specified tolerance*  $\epsilon > 0$

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}^k) < \epsilon$$

- Because  $\ell(\mathbf{w})$  is strongly concave:
  - i.e., because of some technical condition

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \leq \frac{1}{2\lambda} \|\nabla \ell(\mathbf{w})\|_2^2 < \epsilon$$

- Thus, stop when: *I don't know*  $\frac{1}{2\lambda} \|\nabla \ell(\mathbf{w}^k)\|_2^2 < \epsilon$

©Carlos Guestrin 2005-2013

15

## Digression: Logistic regression for more than 2 classes

- Logistic regression in more general case (C classes), where  $Y$  in  $\{1, \dots, C\}$

*for C classes (C-1)(k+1) params*

*forall class  $c \in \{1, \dots, C-1\}$*

$$P(Y=c | \mathbf{x}, \mathbf{w}) \propto e^{w_{c0} + \sum_{i=1}^K w_{ci} x_i}$$

$$P(Y=C | \mathbf{x}, \mathbf{w}) = 1 - \sum_{c=1}^{C-1} P(Y=c | \mathbf{x}, \mathbf{w})$$

*C=2*

*# params to learn: K+1*

$$P(Y=1 | \mathbf{x}, \mathbf{w}) = \frac{e^{w_{10} + \sum_{i=1}^K w_{1i} x_i}}{1 + e^{w_{10} + \sum_{i=1}^K w_{1i} x_i}}$$

*normalizer probs add up to 2*

$$P(Y=0 | \mathbf{x}, \mathbf{w}) = 1 - P(Y=1 | \mathbf{x}, \mathbf{w})$$

$$= \frac{1}{1 + e^{w_{10} + \sum_{i=1}^K w_{1i} x_i}}$$

©Carlos Guestrin 2005-2013

16



## Digression: Logistic regression more generally

$$w = \begin{bmatrix} w_0 & \dots & w_n \end{bmatrix}$$

- Logistic regression in more general case, where  $Y \in \{1, \dots, C\}$

for  $c < C$

$$P(Y = c | \mathbf{x}, \mathbf{w}) = \frac{\exp(w_{c0} + \sum_{i=1}^k w_{ci}x_i)}{1 + \sum_{c'=1}^{C-1} \exp(w_{c'0} + \sum_{i=1}^k w_{c'i}x_i)}$$

← normalize.

for  $c=C$  (normalization, so no weights for this class)

$$P(Y = C | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \sum_{c'=1}^{C-1} \exp(w_{c'0} + \sum_{i=1}^k w_{c'i}x_i)}$$

**Learning procedure is basically the same as what we derived!** Slightly longer derivative

©Carlos Guestrin 2005-2013

17