

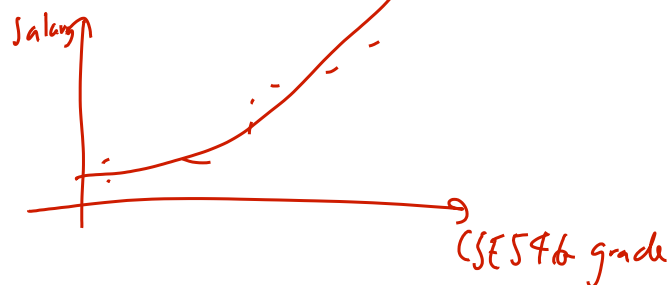
Classification Logistic Regression

Machine Learning – CSE546
Carlos Guestrin
University of Washington

October 9, 2014

©Carlos Guestrin 2005-2014

1

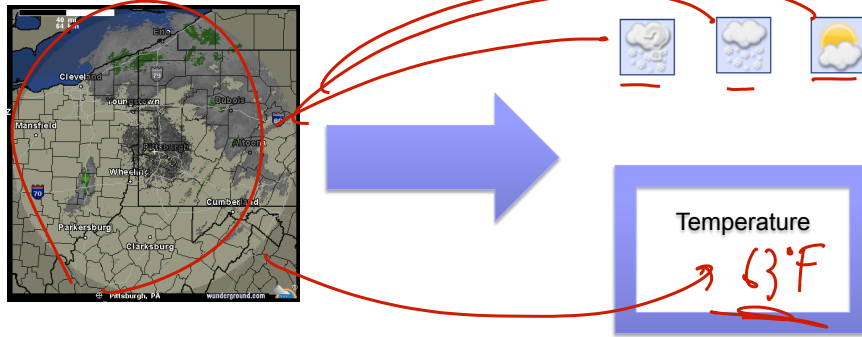


**THUS FAR, REGRESSION:
PREDICT A CONTINUOUS
VALUE GIVEN SOME INPUTS**

©Carlos Guestrin 2005-2014

2

Weather prediction revisited



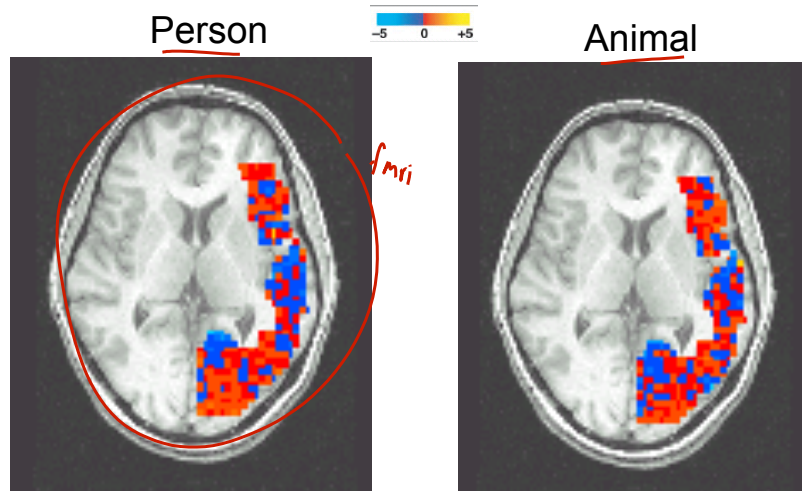
©Carlos Guestrin 2005-2014

3

Reading Your Brain, Simple Example

[Mitchell et al.]

Pairwise classification accuracy: 85%



©Carlos Guestrin 2005-2014

4

Classification

$$X = \{GPA, MC\ grade, \dots\}$$

In reg: $Y = Salary$

in classification: $Y = \{hired, not\ hired\}$

- Learn: $h: X \mapsto Y$

- X – features
- Y – target classes

- Conditional probability: $P(Y|X)$

$$P(Y = hired | X = (GPA = 3.6, MC\ Grade = 3.9))$$

- Suppose you know $P(Y|X)$ exactly, how should you classify?

- Bayes optimal classifier:

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(Y=y | X=x)$$

$$P(hired | 3.6, 3.9) = 0.8$$

$$P(y = not\ hired | 3.6, 3.9) = 0.2$$

$$\Rightarrow g = hired !!$$

- How do we estimate $P(Y|X)$?

©Carlos Guestrin 2005-2014

5

Link Functions

$$X = (x_1, \dots, x_j, \dots)$$

$$X \rightarrow [0, 1]$$

- Estimating $P(Y|X)$: Why not use standard linear regression?

$$P(Y|X) = w_0 + \sum_i w_i x_i$$

$\in [0, 1]$ $\mathbb{R} \in (-\infty, +\infty)$

- Combining regression and probability?

- Need a mapping from real values to $[0, 1]$
- A link function! $g: \mathbb{R} \rightarrow [0, 1]$

Many g options for, but today a simple one

©Carlos Guestrin 2005-2014

6

Logistic Regression

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

Learn $P(Y|X)$ directly

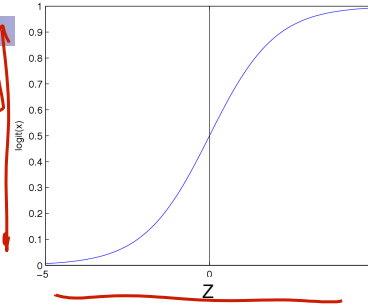
- Assume a particular functional form for link function
- Sigmoid applied to a linear function of the input features: $[0,1]$

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

$$P(Y = 1|X, W) = 1 - P(Y = 0|X, W) = \frac{e^{w_0 + \sum_i w_i X_i}}{1 + e^{w_0 + \sum_i w_i X_i}} = h_i(x)$$

X_i indicators: $X_i = \mathbb{1}(\text{county} = \text{USA})$

Features can be discrete or continuous!



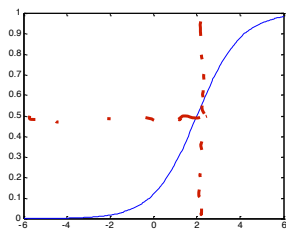
©Carlos Guestrin 2005-2014

7

Understanding the sigmoid

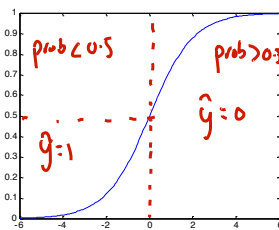
$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

$w_0 = -2, w_1 = -1$



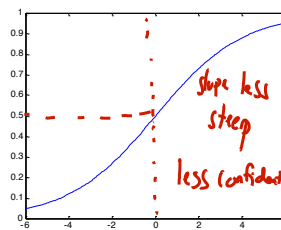
$w_0 + \sum_i w_i x_i = 2$

$w_0 = 0, w_1 = -1$



$w_0 + \sum_i w_i x_i$

$w_0 = 0, w_1 = -0.5$

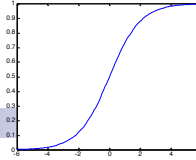


©Carlos Guestrin 2005-2014

8

Logistic Regression – a Linear classifier

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



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

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

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

A red line with a positive slope is drawn across the page, labeled $w_0 + \sum_i w_i x_i$. To the left of the line are '+' signs, and to the right are '-' signs.

©Carlos Guestrin 2005-2014

9

Very convenient! $f > 1 \Rightarrow \ln f > 0$

$$P(Y = 0 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

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

implies

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

implies \log odds ratio

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

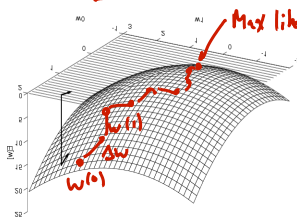
linear classification rule!

©Carlos Guestrin 2005-2014

10

Optimizing concave function – Gradient ascent

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



Gradient: $\nabla_w l(\mathbf{w}) = \left[\frac{\partial l(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial l(\mathbf{w})}{\partial w_n} \right]^T$

Step size, $\eta > 0$

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

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

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

Newton, LBFGS, ...

for HLL

(choice of η ? from theory)
 $\eta = \frac{\alpha}{\sqrt{t}}$
 or
 $\eta = \frac{\alpha}{t}$
 or
 $\eta = \alpha$

Loss function: Conditional Likelihood

- Have a bunch of iid data of the form:

$x \rightarrow$ MC Grade $y \rightarrow$ hired?

$$(x^j, y^j)_{j=1:N} \equiv \mathcal{D} = (\mathcal{D}_X, \mathcal{D}_Y)$$

- Discriminative (logistic regression) loss function:

Conditional Data Likelihood

$$\arg \max_w P(\mathcal{D}_Y | \mathcal{D}_X, w) \stackrel{P(y|x)}{=} \arg \max_w \prod_{i=1}^N P(y^i | x^i, w)$$

$$= \arg \max_w \ln \prod_{i=1}^N P(y^i | x^i, w)$$

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

Expressing Conditional Log Likelihood

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

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

$$l(\mathbf{w}) \equiv \sum_{j=1}^N \ln P(y^j | \mathbf{x}^j, \mathbf{w})$$

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

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

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

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

the best we want to maximize WRT \mathbf{w}

©Carlos Guestrin 2005-2014

13

Maximizing Conditional Log Likelihood

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

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

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

want max \mathbf{w}

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



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

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

Good news: concave functions easy to optimize

gradient ascent

©Carlos Guestrin 2005-2014

14

Maximize Conditional Log Likelihood:

$\frac{d}{dw} \ln f(w) = \frac{f'(w)}{f(w)}$ Gradient ascent $\frac{d}{dw} e^f = f' e^f$

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

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

$P(Y=1 | X^j, w)$

$$\frac{\partial l}{\partial w_k} = \sum_{j=1}^N x_k^j (y^j - P(Y=1 | X^j, w))$$

diff truth & prediction

weighted
by how much feature k plays a role in point j

©Carlos Guestrin 2005-2014

15

Gradient Ascent for LR

initialize $w^{(0)} = 0$ or something else \rightarrow all converge to same place, concavity

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

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

↑ small step ↓ gradient direction

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

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

repeat

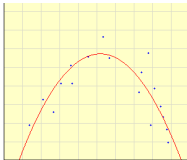
©Carlos Guestrin 2005-2014

16

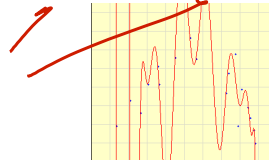
Regularization in linear regression

- 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 $\lambda > 0$:

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

$\|w\|_2^2$
 $\|w\|_1$

©Carlos Guestrin 2005-2014

17

Linear Separability

$\exists w \forall x^i \quad w_0 + \sum_i w_i x_i > 0 \quad g^i = 1$
 $w_0 + \sum_i w_i x_i < 0 \quad g^i = 0$

$w_0 + \sum_i w_i x_i > 0$
 $2w_0 + \sum_i 2w_i x_i > 0$

$w_0 + \sum_i w_i x_i > 0$
 $2w_0 + \sum_i 2w_i x_i < 0$

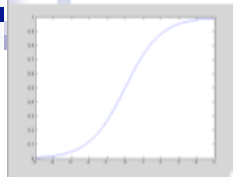
$w_0 + \sum_i w_i x_i > 0$
 $2w_0 + \sum_i 2w_i x_i < 0$

more confident
log likelihood will be higher

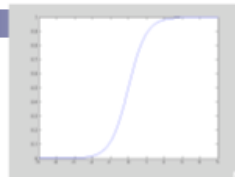
©Carlos Guestrin 2005-2014

18

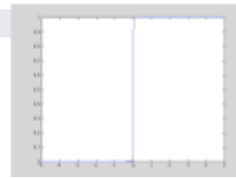
Large parameters → Overfitting



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



$$\frac{1}{1 + e^{-2x}}$$



$$\frac{1}{1 + e^{-100x}}$$

- If data is linearly separable, weights go to infinity

$$P(y=0 | \mathbf{w}, \mathbf{x}^j) \rightarrow 1 \quad \|\mathbf{w}\| \rightarrow \infty$$

- In general, leads to overfitting:
- Penalizing high weights can prevent overfitting...

©Carlos Guestrin 2005-2014

19

Regularized Conditional Log Likelihood

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

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

$\sum_{i=1}^K w_i^2$

- Practical note about w_0 :

don't regularize

- Gradient of regularized likelihood:

$$\frac{\partial \ell}{\partial w_k} = \text{same as before} + \frac{\partial}{\partial w_k}$$

$= -\lambda w_k$

©Carlos Guestrin 2005-2014

20

Standard v. Regularized Updates

- Maximum conditional likelihood estimate

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

Standard update
 $\forall w_i, i > 0$

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

- Regularized maximum conditional likelihood estimate

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

regularization contribution to derivative

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

push w_i towards 0

©Carlos Guestrin 2005-2014

21

Please Stop!! Stopping criterion

$\ell(\mathbf{w}^{(t+1)}) - \ell(\mathbf{w}^{(t)})$

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

- When do we stop doing gradient descent? $\epsilon > 0$

$$\overset{\text{opt}}{\ell(\mathbf{w}^*)} - \ell(\mathbf{w}^{(t)}) < \epsilon$$

- Because $\ell(\mathbf{w})$ is strongly concave:

- i.e., because of some technical condition

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

Don't know this

- Thus, stop when:

$$\|\nabla \ell(\mathbf{w})\|_2^2 < 2\lambda \epsilon$$

implies



$$\eta = \frac{\alpha}{\sqrt{k}}$$

Learning OK

$$\eta = \alpha$$

©Carlos Guestrin 2005-2014

22

Digression: Logistic regression for more than 2 classes

- Logistic regression in more general case (C classes), where Y in $\{0, \dots, C-1\}$

For C class, need $(C-1)(k+1)$ params

If classes $c=1, \dots, C-1$

$$P(Y=c|X, \mathbf{w}) = \frac{e^{w_{c0} + \sum_i w_{ci} x_i}}{1 + \text{something}}$$

For $C=0$

$$P(Y=0|X, \mathbf{w}) = 1 - \sum_{c=1}^{C-1} P(Y=c|X, \mathbf{w})$$

$C=2$, k features

\Rightarrow $k+1$ params to learn w_0, \dots, w_k

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

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

©Carlos Guestrin 2005-2014

23

Digression: Logistic regression more generally

- Logistic regression in more general case, where Y in $\{0, \dots, C-1\}$

for $c > 0$

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

for $c=0$ (normalization, so no weights for this class)

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

Learning procedure is basically the same as what we derived! *Derivative is a little more complicated*

©Carlos Guestrin 2005-2014

24