

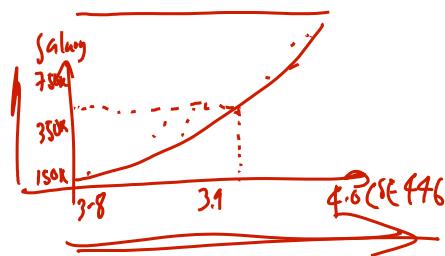
Classification Logistic Regression

Machine Learning – CSE546
Carlos Guestrin
University of Washington

October 9, 2013

©Carlos Guestrin 2005-2013

1

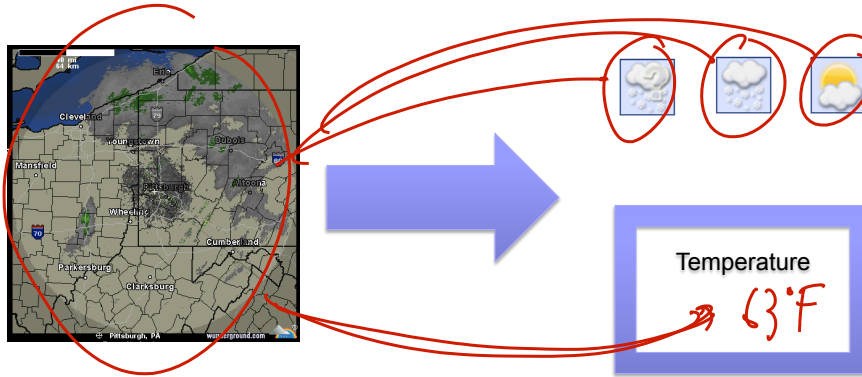


**THUS FAR, REGRESSION:
PREDICT A CONTINUOUS
VALUE GIVEN SOME INPUTS**

©Carlos Guestrin 2005-2013

2

Weather prediction revisited



©Carlos Guestrin 2005-2013

3

Reading Your Brain, Simple Example

[Mitchell et al.]

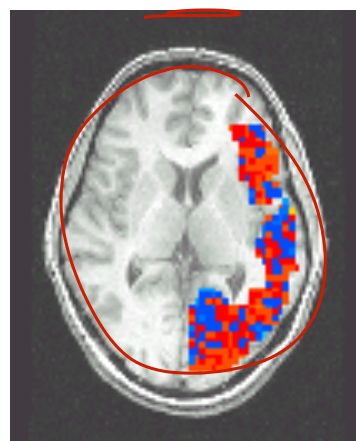
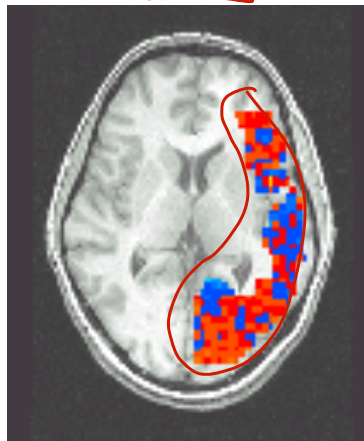
Pairwise classification accuracy: 85%

$C = \{ \text{Person, Animal} \}$

Person



Animal



©Carlos Guestrin 2005-2009

4

Classification

$X \equiv (\text{GPA, ML grade})$

in reg: $Y \equiv \text{Salary}$

in cls: $Y \equiv (\text{hired, not hired})$

Learn: $h: X \mapsto Y$

- X – features
- Y – target classes

Conditional probability: $P(Y|X)$

$P(Y = \text{hired} | X = (\text{GPA} = 3.6, \text{ML grade} = 3.9))$

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

- Bayes optimal classifier:

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

$$P(\text{hired} | 3.6, 3.9) = 0.8$$

$$P(\text{not hired} | 3.6, 3.9) = 0.2$$

$$\Rightarrow \hat{y} = \text{hired}!!$$

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

©Carlos Guestrin 2005-2013

5

Link Functions

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

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

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

Combing regression and probability?

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

Many options, simple choice today

©Carlos Guestrin 2005-2013

6

Logistic Regression

Logistic function (or Sigmoid)

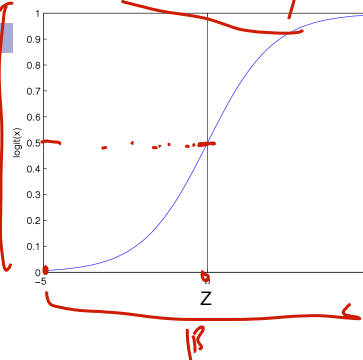
$$\frac{1}{1 + \exp(-z)}$$

Learn $P(Y|X)$ directly

- Assume a particular functional form for link function $(0,1)$
- Sigmoid applied to a linear function of the input features: *choice arbitrary*

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

\mathbb{R}



x_i
Features can be discrete or continuous!

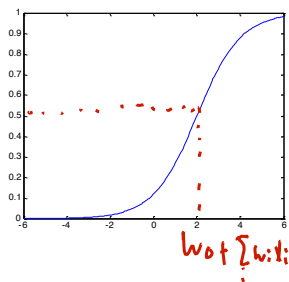
©Carlos Guestrin 2005-2013

7

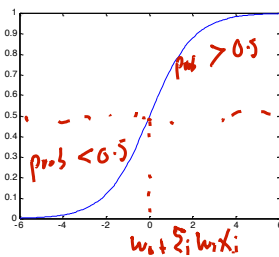
Understanding the sigmoid

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

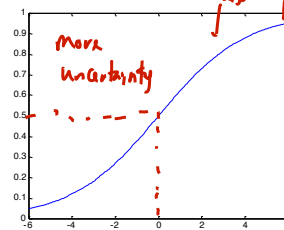
Shift
 $w_0 = -2, w_1 = -1$



$w_0 = 0, w_1 = -1$



$w_0 = 0, w_1 = -0.5$
less steep

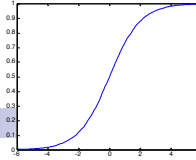


©Carlos Guestrin 2005-2013

8

Logistic Regression – a Linear classifier

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



linear fn separates
positives from negatives

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

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

linear function

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

$w_0 + \sum_i w_i x_i < 0$

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

©Carlos Guestrin 2005-2013

9

Very convenient!

$f > 1 \Rightarrow \ln f > 0$
 called a
 ✓ log linear model

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

implies $P(Y=1|X,w) = 1 - P(Y=0|X,w)$

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

implies

choose $\hat{y}=1$

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

linear
classification
rule!

implies

$0 < \ln$

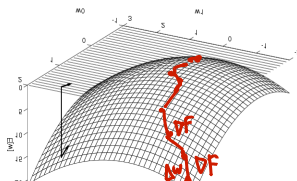
$$\ln \frac{P(Y = 1 | X)}{P(Y = 0 | X)} = w_0 + \sum_i w_i X_i$$

©Carlos Guestrin 2005-2013

10

Optimizing concave function – Gradient ascent *← alternative to coordinate 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

Newtons, LBGs... For your HW, choose a constant η

η choice? in theory, often $\eta = \frac{\alpha}{\epsilon}$ or $\eta = \frac{\alpha}{\sqrt{\epsilon}}$

Loss function: Conditional Likelihood *ln is monotonic*

- Have a bunch of iid data of the form:

X Y
GPA hired?

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

- Discriminative (logistic regression) loss function:

Conditional Data Likelihood

$$\operatorname{argmax}_w P(D_Y | D_X, w) \stackrel{iid}{=} \operatorname{argmax}_w \prod_j P(y^j | x^j, w)$$

$$= \operatorname{argmax}_w \ln \prod_j P(y^j | x^j, w) = \operatorname{argmax}_w \sum_{j=1}^N \ln P(y^j | x^j, w)$$

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

Expressing Conditional Log Likelihood

$$\begin{aligned}
 P(Y=0|X, \mathbf{w}) &= \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)} \\
 P(Y=1|X, \mathbf{w}) &= \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)} \\
 l(\mathbf{w}) &\equiv \sum_{j=1}^N \ln P(y^j | \mathbf{x}^j, \mathbf{w}) \\
 &= \sum_j \begin{cases} \ln P(Y=1 | \mathbf{x}^j, \mathbf{w}) & \text{when } y^j=1 \\ \ln P(Y=0 | \mathbf{x}^j, \mathbf{w}) & \text{when } y^j=0 \end{cases} \\
 \ell(\mathbf{w}) &= \sum_j y^j \ln P(Y=1 | \mathbf{x}^j, \mathbf{w}) + (1 - y^j) \ln P(Y=0 | \mathbf{x}^j, \mathbf{w}) \\
 &= \sum_j y_j \ln \frac{e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}} + (1 - y_j) \ln \frac{1}{1 + e^{w_0 + \sum_i w_i x_i^j}} \\
 &= \sum_{j=1}^N y_j (w_0 + \sum_i w_i x_i^j) - \ln (1 + e^{w_0 + \sum_i w_i x_i^j})
 \end{aligned}$$

want max this wrt w

©Carlos Guestrin 2005-2013

13

Maximizing Conditional Log Likelihood

$$\begin{aligned}
 P(Y=0|X, \mathbf{w}) &= \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)} \\
 P(Y=1|X, \mathbf{w}) &= \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)} \\
 l(\mathbf{w}) &\equiv \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) \\
 &= \sum_j y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))
 \end{aligned}$$

want max w



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

©Carlos Guestrin 2005-2013

14

Maximize Conditional Log Likelihood: Gradient ascent

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

$$\frac{\partial e^f}{\partial w_i} = f' e^f$$

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

$$\frac{\partial l}{\partial w_i} = \sum_{j=1}^N y^j x_i^j - \frac{x_i^j e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}} = \sum_{j=1}^N x_i^j (y^j - P(Y=1 | x^j, w))$$

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

weigh by how much feature i participates in data point x^j

does my prediction match

©Carlos Guestrin 2005-2013

15

Gradient Ascent for LR

$w^{(0)}$ ← initialize, e.g. to β

revisit soon

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

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

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

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

repeat

©Carlos Guestrin 2005-2013

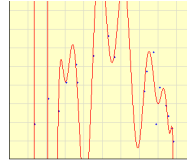
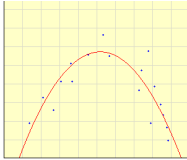
16

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 $\lambda > 0$:

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

$\|w\|_2$
 $\|w\|_1$

©Carlos Guestrin 2005-2013

17

Linear Separability

$\exists w \cdot t_i$
 $w_0 + \sum_i w_i x_i^j > 0$ if $y_j = 1$
 $w_0 + \sum_i w_i x_i^j < 0$ if $y_j = 0$

$w_0 + \sum_i w_i x_i > 0$
 $2w_0 + \sum_i 2w_i x_i > 0$
 $70 \text{ GeV} + w_0 + \sum_i 2w_i x_i > 0 \leftarrow \text{"more sure"}$

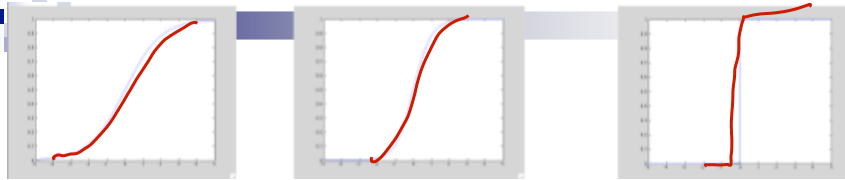
$w_0 + \sum_i w_i x_i < 0$
 $2w_0 + \sum_i 2w_i x_i < 0$

$w_0 + \sum_i w_i x_i$

©Carlos Guestrin 2005-2013

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

$$P(y=0 | \mathbf{w}, x) = \frac{1}{1 + e^{w_0 + \sum w_i x_i}} \rightarrow 1 \text{ for negative examples (and via vs.) as } \|\mathbf{w}\| \rightarrow \infty$$

- In general, leads to overfitting:

- Penalizing high weights can prevent overfitting...



©Carlos Guestrin 2005-2013

19

Regularized Conditional Log Likelihood

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

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

- Practical note about w_0 :

don't regularize w_0 $2w_i$

- Gradient of regularized likelihood:

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

©Carlos Guestrin 2005-2013

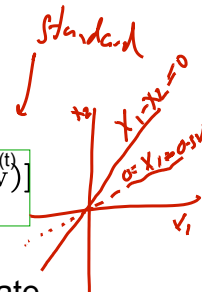
20

Standard v. Regularized Updates

- Maximum conditional likelihood estimate

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

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



- Regularized maximum conditional likelihood estimate

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

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ \underbrace{-\lambda w_i^{(t)}}_{\text{pushes to } 0} + \underbrace{\sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]}_{\text{pushes to max likelihood}} \right\}$$

©Carlos Guestrin 2005-2013

21

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? $\epsilon > 0$

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

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

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

- Thus, stop when:

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

©Carlos Guestrin 2005-2013

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 need

for classes $c \in \{1, \dots, C-1\}$

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

For $c=0$

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

$C=2$

of parameters = $k+1$
for 1 class

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

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

©Carlos Guestrin 2005-2013

23

Digression: Logistic regression more generally

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

for $c > 0$

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

normalization

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

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

just like in 2 class case

Learning procedure is basically the same as what we derived! *derivative a little fancier*

©Carlos Guestrin 2005-2013

24