

Module 5: Classification

Basic Concepts: Risk and Measures of Predictive Accuracy

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 20th, 2014

©Emily Fox 2014

1

The Optimal Prediction

- Assume we *know* the data-generating mechanism $Y|X$
- If our task is prediction, which summary of the distribution $Y|x$ should we report?
For x , what $f(x)$ should we choose to predict Y if we can choose any $f(\cdot)$
- Taking a decision-theoretic framework, consider the **expected loss** predictions are penalized by $L(\cdot, \cdot)$

$$E_{Y|X}[L(Y, f(X))] = E_X\{E_{Y|X}[L(Y, f(X))|X=x]\}$$

truth pred.

— $\hat{f}(\cdot)$ should min \rightarrow

— can min. pointwise

©Emily Fox 2014

2

Continuous Responses

- Expected loss $E_X \{E_{Y|X} [L(Y, f(x)) \mid X = x]\}$

- Example: L_2 $L(Y, f(x)) = (Y - f(x))^2$

Solution: $\hat{f}(x) = E[Y|X]$

focus so far

Proofs:
HW

- Example: L_1 $L(Y, f(x)) = |Y - f(x)|$

Solution: $\hat{f}(x) = \text{median}(Y|x)$

- More generally: L_p $L(Y, f(x)) = \left\{ \int |Y - f(x)|^p \right\}^{1/p}$

©Emily Fox 2014

3

Categorical Responses

$L(k, i)$

- Expected loss $E_X \{E_{Y|X} [L(Y, g(x)) \mid X = x]\}$

- Response: $Y \in \{G_1, \dots, G_K\}$ # of classes

- Same setup, but need new loss function

- Can always represent loss function with $K \times K$ matrix

$$L_{jk} \equiv L(j, k) = \begin{cases} 0, & j=k \\ z_0, & j \neq k \end{cases}$$

- L is zeros on the diagonal and non-negative elsewhere

- Typical loss function:

$$L_{jk} \equiv L(j, k) = \begin{cases} 0, & j=k \\ 1, & j \neq k \end{cases}$$

unit cost for all pass. mistakes

©Emily Fox 2014

4

Optimal Prediction

$$E(X) = \sum X P(X=x)$$

- Expected loss

$$E_X \{ E_{Y|X} [L(Y, g(x)) | X = x] \} = E_X \left(\sum_{k=1}^K L(G_k, g(x)) \cdot P(G_k | X=x) \right)$$

- Again, can minimize pointwise

$$\hat{g}(x) = \underset{g}{\operatorname{argmin}} \sum_{k=1}^K L(g_k, g) P(G_k | X=x)$$

- Example: $K=2$

$$\begin{aligned} & \text{for } g=0: L(1,0)P(G_1|x) + L(0,0)P(G_0|x) > L(1,1)P(G_1|x) + L(0,1)P(G_0|x) \\ & \text{for } g=1: L(1,1)P(G_1|x) + L(0,1)P(G_0|x) > L(1,0)P(G_1|x) + L(0,0)P(G_0|x) \\ & \text{Conclusion: } \hat{g}(x) = 1 \end{aligned}$$

©Emily Fox 2014

5

Optimal Prediction

$$\hat{g}(x) = \underset{g}{\operatorname{argmin}} \sum_{k=1}^K L(G_k, g) \Pr(G_k | X = x)$$

- With 0-1 loss, we straightforwardly get the **Bayes classifier**

$$\hat{g}(x) = \underset{g}{\operatorname{argmin}} [1 - P(g|X=x)] \quad (\text{general})$$

OR

$$\hat{g}(x) = G_k \text{ if } P(G_k | X=x) = \max_g P(g | X=x)$$

(classify to most probable class)

©Emily Fox 2014

6

Optimal Prediction

$$\hat{g}(x) = \mathcal{G}_k \quad \text{if} \quad \Pr(\mathcal{G}_k | X = x) = \max_g \Pr(g | X = x)$$

■ How to approximate the optimal prediction?

- Don't actually have $p(Y | X = x)$

■ Nearest neighbor approach

- Look at k -nearest neighbors with majority vote to estimate

$$P(G_m | X=x) \approx \frac{1}{K} \sum_{x_i \in \text{nbhd}(x)} \mathbb{I}(y_i = m)$$

classify w/ Largest $P(G_m | X=x)$
(most common label of k -NN)

before (L_2)

$$E[Y | X]$$

$$= \text{avg}(y_i | x_i)$$

$\in \text{nbhd}(x)$

©Emily Fox 2014

7

Optimal Prediction

$$\hat{g}(x) = \mathcal{G}_k \quad \text{if} \quad \Pr(\mathcal{G}_k | X = x) = \max_g \Pr(g | X = x)$$

■ How to approximate the optimal prediction?

- Don't actually have $p(Y | X = x)$

■ Model-based approach

- Introduce indicators for each class:

- Consider squared-error loss: $\hat{f}(X) = E[Y | X]$

$$E[Y_k | X] = P(Y = G_k | X)$$

- Bayes classifier is equivalent to standard regression and L_2 loss, followed by classification to largest fitted value

$$\hat{f}(x) = \begin{bmatrix} 0.1 \\ 0.02 \\ \vdots \end{bmatrix} \rightarrow \text{largest} \rightarrow \text{class}$$

- Works in theory, but not in practice... Will look at many other approaches (e.g., logistic regression)

©Emily Fox 2014

8

Measuring Accuracy of Classifier

- For a given classifier, how do we assess how well it performs?
- For 0-1 loss, the generalization error is

$$E_{x,y} [g(x) \neq y] = P_{x,y} (g(x) \neq y)$$

with empirical estimate

$$\frac{1}{m} \sum_{i=1}^m \mathbb{I}(g(x_i) \neq y_i)$$

chosen classifier

- Consider binary response and some useful summaries

decision rule

$$y = \begin{cases} 0 \\ 1 \end{cases} \quad \begin{matrix} \text{no disease} \\ \text{disease} \end{matrix}$$

$$g(x) = \begin{cases} 0 \\ 1 \end{cases} \quad \begin{matrix} \text{predict no disease} \\ \text{predict disease} \end{matrix}$$

©Emily Fox 2014 9

Measuring Accuracy of Classifier

- Sensitivity:

prob. of pred. disease for a diseased individual

$$P(G(x)=1 | Y=1)$$

- Specificity:

no disease given individual's not diseased

$$P(G(x)=0 | Y=0)$$

- False positive rate:

$$P(g(x)=1 | Y=0)$$

- True positive rate:

$$P(g(x)=1 | Y=1)$$

- Connections:

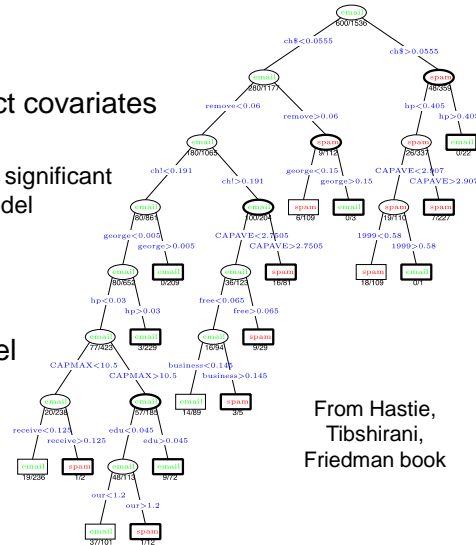
$$\text{sensitivity} = \text{TPR}, \text{ specificity} = 1 - \text{FPR}$$

©Emily Fox 2014 10

Classification Tree Spam Example

- Resulting tree of size 17
- Note that there are 13 distinct covariates split on by the tree
 - 11 of these overlap with the 16 significant predictors from the additive model previously explored
- Overall error rate (9.3%) is higher than for additive model

True	Predicted	
	email	spam
email	57.3%	4.0%
spam	5.3%	33.4%



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

11

Classification Tree Spam Example

- Think of **spam** and **email** as presence and absence of disease

- Using equal losses

Sensitivity = $100 \times \frac{33.4}{33.4 + 5.3} = 86.3\%$

Specificity = $100 \times \frac{57.3}{57.3 + 4.0} = 93.4\%$

True	Predicted	
	email	spam
email	57.3%	4.0%
spam	5.3%	33.4%

From Hastie, Tibshirani, Friedman book

- By varying L_{01} and L_{10} , can increase/decrease sensitivity and decrease/increase specificity of rule

- Which do we want here? avoid marking 'email' as 'spam'

- How? $L_{01} > L_{10} = 1$... high specificity

- Change in rule at leaf:

predict 'spam' if proportion of 'spam' at leaf $\geq \frac{L_{01}}{L_{01} + L_{10}}$

©Emily Fox 2014

12

ROC Curves

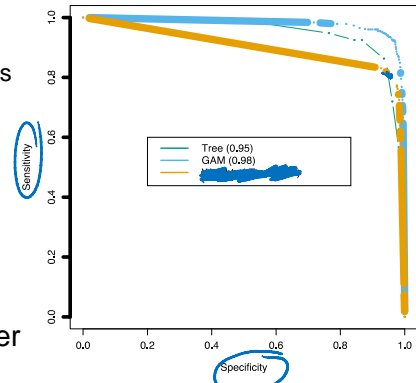
- **Receiver operating characteristic** (ROC) curve summarizes tradeoff between sensitivity and specificity
 - Plot of sensitivity vs. specificity as a function of params of classification rule

- Example: vary L_0 in $[0.1, 10]$
 - Want specificity near 100%, but in this case sensitivity drops to about 50%

- Summary = area under the curve

- Tree = 0.95
- GAM = 0.98

- Instead of Bayes rule at leaf, better to account for unequal losses in constructing tree



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

13

What you need to know

- Again, goal framed as minimizing expected loss
- Loss here is summarized by $K \times K$ matrix L (K classes)
 - Common choice = 0-1 loss
- Bayes classifier chooses most probable class (intuitive)
- Measures of predictive performance:
 - Sensitivity, specificity, true positive rate, false positive rate
 - ROC curve and area under the curve

©Emily Fox 2014

14

Readings

- Wakefield – 10.3.2, 10.4.2, 12.8.4
- Hastie, Tibshirani, Friedman – 9.2.3, 9.2.5, 2.4

Module 5: Classification

Linear Methods: Logistic Regression

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 20th, 2014

Link Functions

Focus on $Y \in \{0,1\}$

or generally

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

$Y \in \{1, \dots, k\}$

$$p(Y|X) = \beta_0 + \sum_j \beta_j h_j(x)$$

$\in [0,1]$

range: $(-\infty, \infty)$

BAD

- Combining regression and probability?

- Need a mapping from real values to $[0,1]$
- A link function!

$g: \mathbb{R} \rightarrow [0,1]$
many options, but here's a useful one

©Emily Fox 2014

17

Logistic Regression

Logistic function

(or Sigmoid):

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

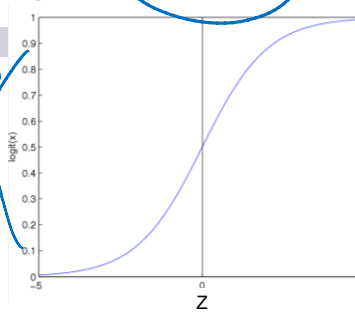
- Learn $p(Y|X)$ directly

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

$$p(y=0 | x, \beta) = \frac{1}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

choice/Label

z : linear $(-\infty, \infty)$ just like in regression \mathbb{R}



$$p(Y=1|X) = \frac{e^{(\dots)}}{1 + e^{(\dots)}}$$

$\beta_0 + \sum_j \beta_j x_j$
not bounded, could be neg.

after applying the sigmoid, output is in $[0,1]$

Covariates can be discrete or continuous!

©Emily Fox 2014

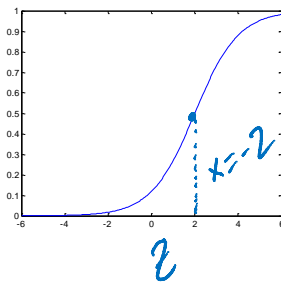
18

Understanding the Sigmoid

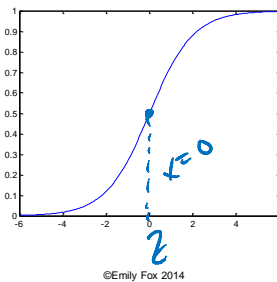
$$g(\beta_0 + \sum_j \beta_j x_j) = \frac{1}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

$\downarrow d=1$

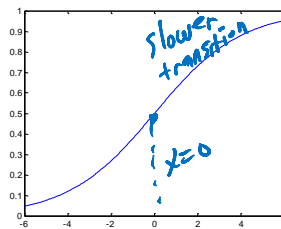
$\beta_0=-2, \beta_1=-1$



$\beta_0=0, \beta_1=-1$



$\beta_0=0, \beta_1=-0.5$

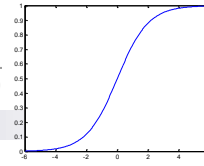


©Emily Fox 2014

19

Logistic Regression – a Linear classifier

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



$$g(\beta_0 + \sum_j \beta_j x_j) = \frac{1}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

$P(y=1|x, \beta)$

$\beta_0 + \sum_j \beta_j x_j = 0$
hyperplane

$\beta_0 + \sum_j \beta_j x_j > 0$
 $\rightarrow g(\cdot) < 0.5$
 $\rightarrow P(y=1) < 0.5$
 \rightarrow predict class 0

$(\cdot) < 0$
 $\rightarrow g(\cdot) > 0.5$
 \rightarrow predict class 1

©Emily Fox 2014

20

Very convenient!

$$p(y = 0 | x, \beta) = \frac{1}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

implies

$$p(y = 1 | x, \beta) = \frac{\exp(\beta_0 + \sum_j \beta_j x_j)}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

Examine ratio:

$$\frac{p(y = 1 | x, \beta)}{p(y = 0 | x, \beta)} = \exp(\beta_0 + \sum_j \beta_j x_j)$$

implies

$$\log \frac{p(y = 1 | x, \beta)}{p(y = 0 | x, \beta)} = \beta_0 + \sum_j \beta_j x_j$$

linear
classification
rule!

class 1 wins
else class 0
(under 0-1 loss)

class 1
wins!

linear

©Emily Fox 2014

21

Loss Function: Conditional Likelihood

- Have a bunch of iid data of the form:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \text{ iid } i=1, \dots, n$$

$$= (D_X, D_Y)$$

- Discriminative (logistic regression) loss function:
Conditional Data Likelihood

$$\arg\max_{\beta} p(D_Y | D_X, \beta) = \arg\max_{\beta} \log \prod_{i=1}^n p(y_i | x_i, \beta)$$

$$= \arg\max_{\beta} \sum_{i=1}^n \log p(y_i | x_i, \beta)$$

$$\log p(D_Y | D_X, \beta) = \sum_{i=1}^n \log p(y_i | x_i, \beta)$$

©Emily Fox 2014

22

$$\begin{aligned}
 l(\beta) &= \sum_i \log p(y_i | x_i, \beta) \\
 &= \sum_i \begin{cases} \log p(y=1 | x_i, \beta) & \text{if } y_i=1 \\ \log p(y=0 | x_i, \beta) & \text{if } y_i=0 \end{cases} \\
 l(\beta) &= \sum_i y_i \log p(y=1 | x_i, \beta) + (1 - y_i) \log p(y=0 | x_i, \beta) \\
 &= \sum_i y_i \log \left(\frac{\exp(\beta_0 + \sum_j \beta_j x_{ij})}{1 + \exp(\beta_0 + \sum_j \beta_j x_{ij})} \right) + (1 - y_i) \log \left(\frac{1}{1 + \exp(\beta_0 + \sum_j \beta_j x_{ij})} \right) \\
 &= \sum_i y_i (\beta_0 + \sum_j \beta_j x_{ij}) - \log(1 + \exp(\beta_0 + \sum_j \beta_j x_{ij}))
 \end{aligned}$$

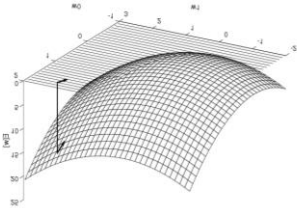
$$l(\beta) = \sum_i \log p(y_i | x_i, \beta) = \sum_i y_i (\beta_0 + \sum_j \beta_j x_{ij}) - \log(1 + \exp(\beta_0 + \sum_j \beta_j x_{ij}))$$

Handwritten notes on the slide:

- Blue circles around β_0 and β_j in the equation.
- Blue arrows pointing from the circles to the text "parameters to opt." with a checkmark.
- Blue arrow pointing from the β_j circle to the text "fixed".

Optimizing Concave Function – Gradient Ascent

- Conditional likelihood for logistic regression is concave
- Find optimum with gradient ascent



$$\text{Gradient: } \nabla_{\beta} l(\beta) = \left[\frac{\partial l(\beta)}{\partial \beta_0}, \dots, \frac{\partial l(\beta)}{\partial \beta_d} \right]'$$

Step size, $\eta > 0$

$$\text{Update rule: } \Delta \beta = \eta \nabla_{\beta} l(\beta)$$

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \frac{\partial l(\beta)}{\partial \beta_j}$$

↑ newest. ↑ oldest.

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

often, η gets smaller w/ iteration
 $\eta = \frac{1}{t} (\text{const.})$

©Emily Fox 2014

25

Maximize Conditional Log Likelihood:

$\log \zeta = \log f(\beta)$ Gradient ascent $\frac{\partial}{\partial \beta} \log f(\beta) = \frac{f'(\beta)}{f(\beta)}$

$$l(\beta) = \sum_i y_i (\beta_0 + \sum_j \beta_j x_{ij}) - \log(1 + \exp(\beta_0 + \sum_j \beta_j x_{ij}))$$

$$\nabla l(\beta): \frac{\partial l(\beta)}{\partial \beta_j} = \sum_{i=1} y_i x_{ij} - \frac{x_{ij} \exp(\beta_0 + \sum_j \beta_j x_{ij})}{1 + \exp(\beta_0 + \sum_j \beta_j x_{ij})}$$

$$= \sum_i x_{ij} (y_i - \hat{p}(y_i | x_{ij}, \beta))$$

↑ weighted by contribution of j th covariate to y_i

how far is my pred. from truth

©Emily Fox 2014

26

Gradient Ascent for LR

start w/ $\beta^0 = 0$

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

$$\beta_0^{(t+1)} \leftarrow \beta_0^{(t)} + \eta \sum_i (y_i - \hat{p}(y = 1 | x_i, \beta^{(t)}))$$

For $j=1, \dots, d$,

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \sum_i x_{ij} (y_i - \hat{p}(y = 1 | x_i, \beta^{(t)}))$$

repeat

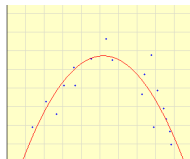
©Emily Fox 2014

27

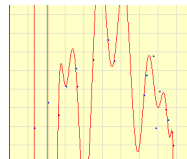
Regularization in Linear Regression

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

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



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



even for
 $n \gg p$,
 p large

- Regularized** or **penalized** regression aims to impose a “complexity” penalty by penalizing large weights

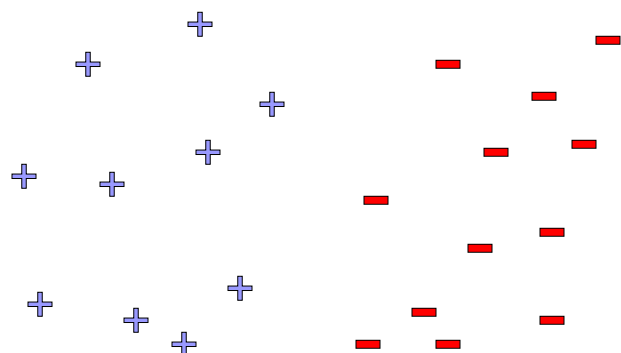
- “Shrinkage” method

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (y_i - (\beta_0 + \beta^T x_i))^2 + \lambda \|\beta\|$$

©Emily Fox 2014

28

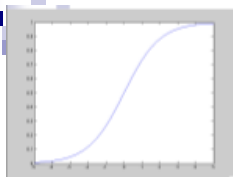
Linear Separability



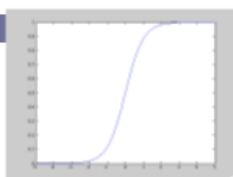
©Emily Fox 2014

29

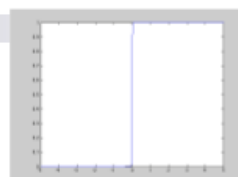
Large Parameters → Overfitting



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



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



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

- If data is linearly separable, weights go to infinity

□ In general, leads to overfitting:

- Penalizing high weights can prevent overfitting...

©Emily Fox 2014

30

Regularized Conditional Log Likelihood

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

$$l(\beta) = \log \prod_{i=1}^n p(y_i | x_i, \beta) - \frac{\lambda}{2} \|\beta\|_2^2$$

- Practical note about β_0 :
- Gradient of regularized likelihood:

©Emily Fox 2014

31

Standard v. Regularized Updates

- Maximum conditional likelihood estimate

$$\hat{\beta} = \arg \max_{\beta} \log \prod_{i=1}^n p(y_i | x_i, \beta)$$

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \sum_i x_{ij} \left(y_i - \hat{p}(y = 1 | x_i, \beta^{(t)}) \right)$$

- Regularized maximum conditional likelihood estimate

$$\hat{\beta} = \arg \max_{\beta} \log \prod_{i=1}^n p(y_i | x_i, \beta) - \frac{\lambda}{2} \sum_{j=1}^d \beta_j^2$$

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + \sum_i x_{ij} \left(y_i - \hat{p}(y = 1 | x_i, \beta^{(t)}) \right) \right\}$$

©Emily Fox 2014

32

Stopping Criterion

$$l(\beta) = \log \prod_{i=1}^n p(y_i | x_i, \beta) - \frac{\lambda}{2} \|\beta\|_2^2$$

- When do we stop doing gradient ascent?

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

- ☐ i.e., because of some technical condition

$$l(\beta^*) - l(\beta) \leq \frac{1}{2\lambda} \|\nabla l(\beta)\|_2^2$$

- Thus, stop when:

©Emily Fox 2014

33

Digression:

Logistic Regression for $K > 2$

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

©Emily Fox 2014

34

Digression: Logistic Regression for $K > 2$

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

for $k < K$

$$p(y = k | \mathbf{x}, \beta) = \frac{\exp(\beta_{k0} + \sum_{j=1}^d \beta_{kj} x_j)}{1 + \sum_{k'=1}^{K-1} \exp(\beta_{k'0} + \sum_{j=1}^d \beta_{k'j} x_j)}$$

for $k=K$ (normalization, so no weights for this class)

$$p(y = K | \mathbf{x}, \beta) = \frac{1}{1 + \sum_{k'=1}^{K-1} \exp(\beta_{k'0} + \sum_{j=1}^d \beta_{k'j} x_j)}$$

**Estimation procedure is basically the same
as what we derived!**

©Emily Fox 2014

35

The Cost, The Cost!!! Think about the cost...

- What's the cost of a gradient update step for LR???

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + \sum_i x_{ij} (y_i - \hat{p}(y = 1 | x_i, \beta^{(t)})) \right\}$$

©Emily Fox 2014

36

Gradient ascent in Terms of Expectations

- “True” objective function:

$$l(\beta) = E_x[l(\beta, x)] = \int p(x)l(\beta, x)dx$$

- Taking the gradient:
- “True” gradient ascent rule:
- How do we estimate expected gradient?

©Emily Fox 2014

37

SGD: Stochastic Gradient Ascent (or Descent)

- “True” gradient: $\nabla l(\beta) = E_x[\nabla l(\beta, x)]$
- Sample based approximation:
- What if we estimate gradient with just one sample???
 - ☐ Unbiased estimate of gradient
 - ☐ Very noisy!
 - ☐ Called stochastic gradient ascent (or descent)
 - Among many other names
 - ☐ VERY useful in practice!!!

©Emily Fox 2014

38

Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

$$E_x[l(\beta, x)] = E_x \left[\log p(y | x, \beta) - \frac{\lambda}{2} \|\beta\|_2^2 \right]$$

- Batch gradient ascent updates:

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + \frac{1}{n} \sum_{i=1}^n x_{ij} \left(y_i - \hat{p}(y = 1 | x_i, \beta^{(t)}) \right) \right\}$$

- Stochastic gradient ascent updates:

- Online setting:

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + x_{i(t),j} \left(y_{i(t)} - \hat{p}(y = 1 | x_{i(t)}, \beta^{(t)}) \right) \right\}$$

©Emily Fox 2014

39

What you should know...

- Classification: predict discrete classes rather than real values
- Logistic regression model: Linear model
 - Logistic function maps real values to [0,1]
- Optimize conditional likelihood
- Gradient computation
- Overfitting
- Regularization
- Regularized optimization
- Cost of gradient step is high, use stochastic gradient descent

©Emily Fox 2014

40