

## Module 5: Classification

# Mixture Models for Classification

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 27<sup>th</sup>, 2014

©Emily Fox 2014

1

## Overview of Classification So Far

- Supervised methods

Generative

Discriminative

LDA, QDA,  
KDE for class,  
Naive Bayes

logistic regr.  
CART

- Objectives:

$$\max P(y, x)$$

$$\max P(y|x)$$

- Unsupervised methods (generative)

mixture models

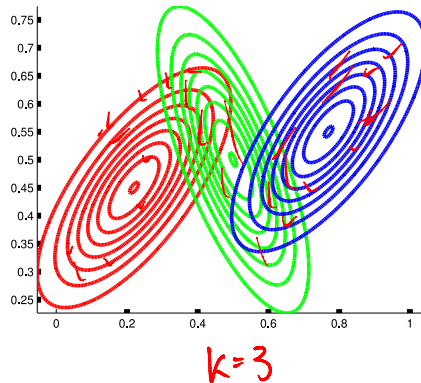
©Emily Fox 2014

2

# Density as Mixture of Gaussians

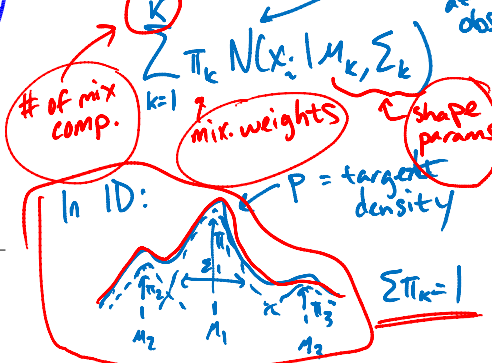
- Approximate density with a mixture of Gaussians

Mixture of 3 Gaussians



$$p(x_i | \pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k)$$

Gauss. kernel, just like in KDE, but not centered at obs.



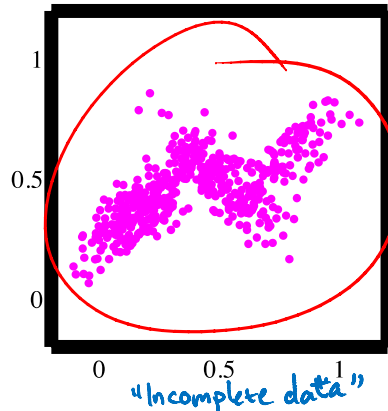
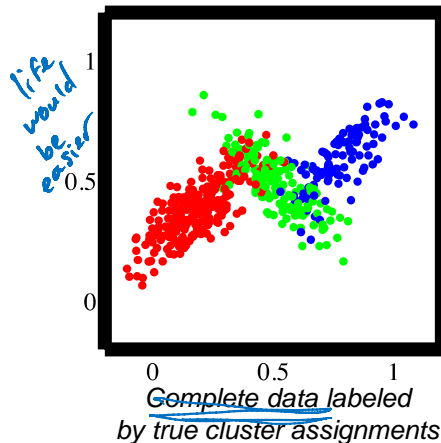
©Emily Fox 2014

3

# Clustering our Observations

- Imagine we have an assignment of each  $x_i$  to a Gaussian

Our actual observations

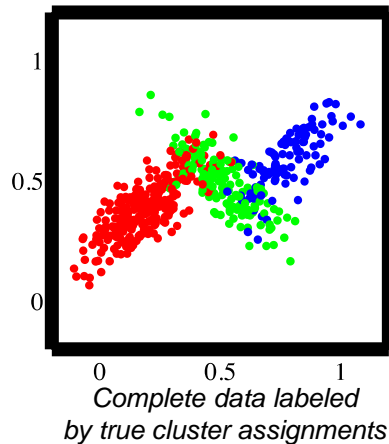


C. Bishop, Pattern Recognition & Machine Learning

©Emily Fox 2014

# Clustering our Observations

- Imagine we have an assignment of each  $x_i$  to a Gaussian



- Introduce latent cluster indicator variable  $z_i$

$$z_i \in \{1, \dots, K\}$$

$$Pr(z_i = k) = \pi_k$$

- Then we have

$$p(x_i | z_i, \pi, \mu, \Sigma) =$$

$$N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

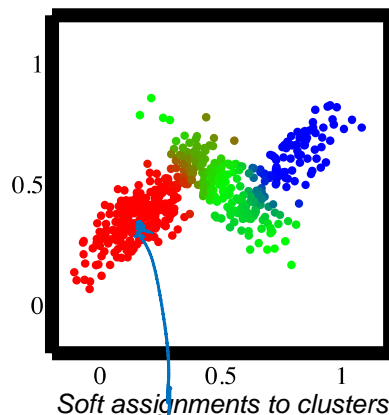
param. est. is easy if we have  $\{z_i\}$   
 $\Rightarrow$  decoupled into  $K$  Gauss. est.

think of as class indicator  $y_i$ , but no training data labels

C. Bishop, Pattern Recognition & Machine Learning

# Clustering our Observations

- We must infer the cluster assignments from the observations



- Posterior probabilities of assignments to each cluster \*given\* model parameters:

$$r_{ik} = p(z_i = k | x_i, \pi, \theta) =$$

$$= \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_i | \mu_j, \Sigma_j)}$$

motivates an iterative alg.

"responsibilities"

(soft cluster prob.)

C. Bishop, Pattern Recognition & Machine Learning

# Mixture Models for Classification

- Can use mixture models as a generative classifier in the unsupervised setting

- EM algorithm = iteratively:

- Estimate responsibilities given parameter estimates

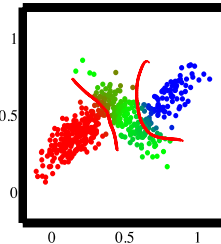
$$\hat{r}_{ik} = \frac{\hat{\pi}_k N(x_i, \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{\ell} \hat{\pi}_{\ell} N(x_i, \hat{\mu}_{\ell}, \hat{\Sigma}_{\ell})}$$

- Maximize parameters given responsibilities

- For classification, threshold the estimated responsibilities

- E.g.,  $\hat{g}(x_i) = \arg \max_k \hat{r}_{ik}$

- Note: allows non-linear boundaries as in QDA

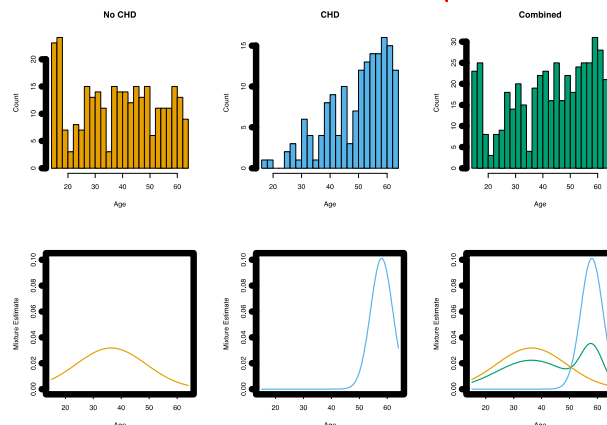


©Emily Fox 2014

7

## Example: Heart Disease Data

- Binary response = CHD (coronary heart disease)
- Predictor = ~~systolic blood pressure~~ *age*



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

8



# What you need to know

- Discriminative vs. Generative classifiers
- LDA and QDA assume Gaussian class-conditional densities
  - Results in linear and quadratic decision boundaries, respectively
- KDE for classification
  - Challenging in areas with little data or in high dimensions
  - Estimating class-conditionals is not optimizing classification objective
- Naïve Bayes assumes factored form
  - Results in log odds that have GAM form
- Mixture models allow for unsupervised generative approach

$$P(X|Y=k)$$

$$(\Sigma_k = \Sigma \text{ for LDA or not})$$

$$P(X|Y=k) = \prod_{i=1}^d P(x_i|Y=k)$$

©Emily Fox 2014

9

# Readings

- Hastie, Tibshirani, Friedman – 4.3, 4.4.5, 6.6.2-6.6.3, 6.8

©Emily Fox 2014

10

## Module 5: Classification

# Online Learning Perceptron Algorithm

STAT/BIOSTAT 527, University of Washington

Emily Fox

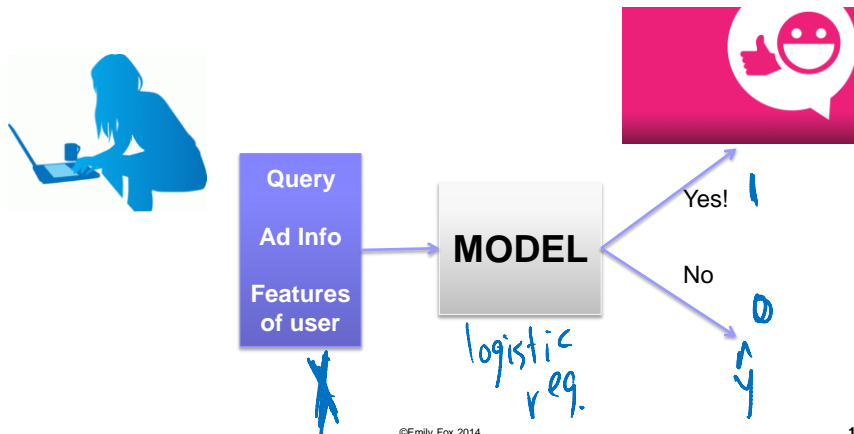
May 27<sup>th</sup>, 2014

©Emily Fox 2014

11

## Estimating Click Probabilities

- **Goal:** Predict whether a person clicks on an ad
- **Basic approach:** Logistic regression



©Emily Fox 2014

12

## Challenge 1: Complexity of Computing Gradients *in terms of n, d*

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

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

*for each j* (stochastic gradient ascent) *O(d)* *O(nd)*

Naively,  $O(nd^2)$

but cache  $\hat{p}$  (same  $\forall j$ )  $\rightarrow O(nd)$

However, if  $n$  is huge (or streaming), this is very slow (infeasible) per little gradient step

©Emily Fox 2014

13

## Challenge 2: Data is streaming

- Assumption thus far: **Batch data** = *have dataset, run model*

- But, e.g., click prediction for ads is a **streaming** data task:

- User enters query, and ad must be selected:

- Observe  $x_i$ , and must predict  $y_i$



*webpage*

*features*

*predicted class  $\hat{y}_i$*

*then observe  $y_i$*

- User either clicks or doesn't click on ad:

- Label  $y_i$  is revealed afterwards

- Google gets a reward if user clicks on ad

*reward is just like 0-1 loss in class setting*

*goal: which ads have a high prob. of click*

- Weights must be updated for next time:

*$\beta^{(t+1)} = \beta^{(t)} + \Delta$  ?? what is it?*

©Emily Fox 2014

14

# Online Learning Problem

- At each time step  $t$ :

- Observe features (covariates) of data point:

- Note: many assumptions are possible, e.g., data is iid, data is adversarially chosen... details beyond scope of course

- Make a prediction:

- Note: many models are possible, we focus on linear models

$$\beta_0^{(t)} + \sum \beta_j^{(t)} x_{tj} > 0 \rightarrow \text{"click"} \quad \parallel \quad \beta^{(t)} \cdot x_t > 0$$

*vec. vec.*

- Observe true label:

- Note: other observation models are possible, e.g., we don't observe the label directly, but only a noisy version... Details beyond scope of course

$y_t \rightarrow 1$  clicked  
 $y_t \rightarrow 0$  not clicked

- Update model:

$$\beta^{(t+1)} = \beta^{(t)} + \Delta^{(t)} \quad \text{WHAT?}$$

©Emily Fox 2014

15

# The Perceptron Algorithm [Rosenblatt '58, '62]

- Classification setting:  $y$  in  $\{-1, +1\}$

*(diff. from  $\{0, 1\}$ , practical change)*

- Linear model

- Prediction:

$$\hat{y} = \text{sign}(\beta \cdot x)$$

- Training:

- Initialize weight vector:

$$\beta^{(0)} = 0 \quad (\text{or smarter})$$

- At each time step:

- Observe covariates:

$$x_t \quad \hat{y} = \text{sign}(\beta^{(t)} \cdot x)$$

- Make prediction:

- Observe true class:

$$y_t \leftarrow \text{true label}$$

- Update model:

- If prediction is not equal to truth

$$\text{if } \hat{y} \neq y_t, \beta^{(t+1)} = \beta^{(t)} \\ \text{else, } \beta^{(t+1)} = \beta^{(t)} + (y_t x_t)$$

*i.e. making a mistake*

©Emily Fox 2014

16

## Intuition

If  $\hat{y} = y_t$ ,

$$\beta^{(t+1)} \leftarrow \beta^{(t)}$$

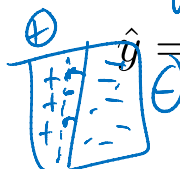
else

$$\beta^{(t+1)} \leftarrow \beta^{(t)} + y_t x_t$$

by adding " $x_t$ " to  $\beta^{(t)}$ , we ↑  $\beta^{(t)} \cdot x_t$  the most.

(similar logic when  $y_t = -1$ )

$$\hat{y} = \text{sign}(\beta^{(t)} \cdot x_t)$$



- Why is this a reasonable update rule?

If mistake when  $y_t = 1, \hat{y} = -1 \rightarrow \beta^{(t)} \cdot x_t < 0$  ... want  $\beta^{(t)} \cdot x_t > 0$

(what " $u$ " maximizes  $u \cdot x_t$ ?)

unit vector ( $u = x_t$ )

$\cos(\theta(u, x_t)) \propto u \cdot x_t$

©Emily Fox 2014

17

## Which weight vector to report?

- Practical problem for all online learning methods
- Suppose you run online learning method and want to sell your learned weight vector... Which one do you sell???

- Last one?

$\beta^{(t)}$ ? no, noisy b/c influenced by last mistake.

- Random

$\beta^{(rand)}$  no

- Average

$$\hat{\beta} = \frac{1}{T} \sum_{t=1}^T \beta^{(t)}$$

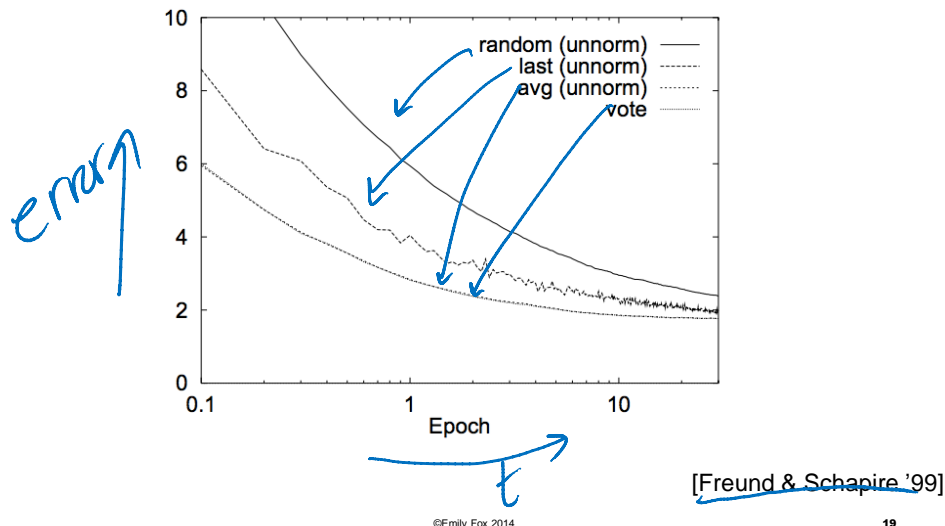
(easy to maintain)

- Voting + more advanced

©Emily Fox 2014

18

## Choice can make a huge difference!!



## Mistake Bounds Why does it work?

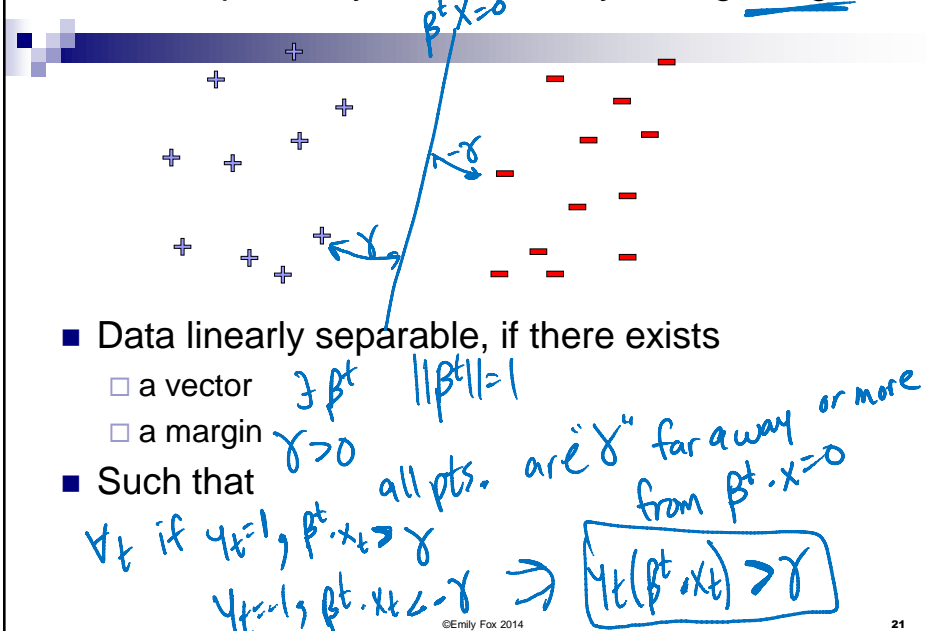
- Algorithm “pays” every time it makes a mistake:

Loss fcn in online setting:  
# of mistakes up to time  $T$

→ Google “pays” for its mistakes

- How many mistakes is it going to make? in lifetime  
“mistake bound”

## Linear Separability: More formally, Using Margin



## Perceptron Analysis: Linearly Separable Case

### ■ Theorem [Block, Novikoff]:

- Given a sequence of labeled examples:  $(x_1, y_1), \dots, (x_n, y_n)$

- Each covariate vector has bounded norm:

- If dataset is linearly separable:

- $\exists \beta^* \text{ s.t. } \|\beta^*\|=1 \text{ and } y_t(\beta^* \cdot x_t) > \gamma, \gamma > 0.$
- Then the number of mistakes made by the online perceptron on this sequence is bounded by

$\left(\frac{R}{\gamma}\right)^2$  **FOREVER**  
 crazy!  
 constant... doesn't depend on T!!

## Perceptron Proof for Linearly Separable case

- Every time we make a mistake, we get  $\gamma$  closer to  $\beta^*$ :

- Mistake at time  $t$ :  $\beta^{(t+1)} = \beta^{(t)} + y_t x_t$

- Taking dot product with  $\beta^*$ :

- Thus after  $m$  mistakes:

- Similarly, norm of  $\beta^{(t+1)}$  doesn't grow too fast:

- $\|\beta^{(t+1)}\|^2 = \|\beta^{(t)}\|^2 + 2y_t(\beta^{(t)} \cdot x_t) + \|x_t\|^2 \leq R^2$

- Thus, after  $m$  mistakes:

- Putting all together:

by induction,  $\beta^* \cdot \beta^{(t+1)} \geq m\gamma$

$m\gamma \leq \beta^* \cdot \beta^{(m)} \leq \|\beta^*\| \|\beta^{(m)}\| \leq \sqrt{m} R$

$\rightarrow m\gamma \leq \sqrt{m} R \rightarrow m \leq \left(\frac{R}{\gamma}\right)^2$

©Emily Fox 2014

23

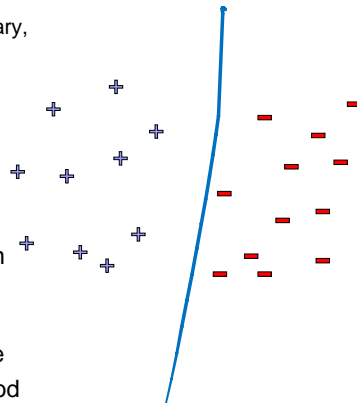
## Beyond Linearly Separable Case

- Perceptron algorithm is super cool!

- No assumption about data distribution!
    - Could be generated by an oblivious adversary, no need to be iid
  - Makes a fixed number of mistakes, and it's done for ever!
    - Even if you see infinite data

- However, real world not linearly separable

- Can't expect never to make mistakes again
  - Analysis extends to non-linearly separable case
  - Very similar bound, see Freund & Schapire
  - Converges, but ultimately may not give good accuracy (make many many many mistakes)



©Emily Fox 2014

24



# What is the Perceptron Doing???

- When we discussed logistic regression:
  - Started from maximizing conditional log-likelihood

$$\max_{\beta} P(y|x)$$

- When we discussed the perceptron:
  - Started from description of an algorithm

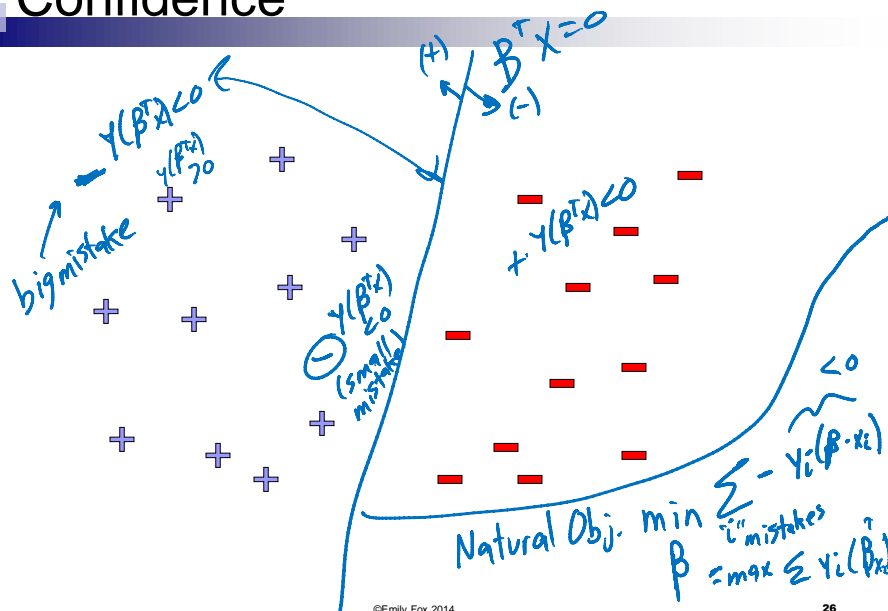
- What is the perceptron optimizing????

(loss fcn.s.)

©Emily Fox 2014

25

## Perceptron Prediction: Margin of Confidence

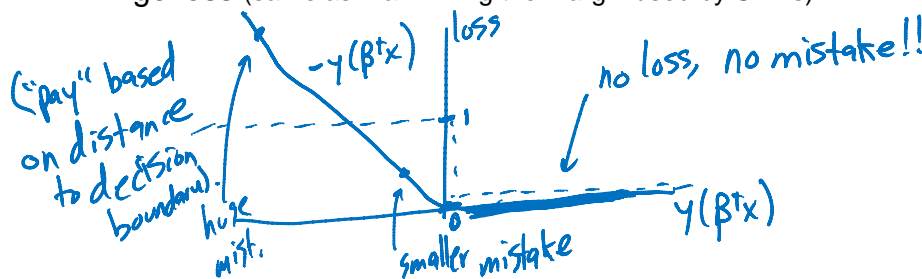


©Emily Fox 2014

26

# Hinge Loss

- Perceptron prediction:  $\text{sign}(\beta^T x)$
- Makes a mistake when:  $y(\beta^T x) < 0 \rightarrow L(\beta, x) = \begin{cases} 0 & y(\beta^T x) \geq 0 \\ -y(\beta^T x) & \text{"other"} \rightarrow (-y(\beta^T x)) \end{cases}$
- Hinge loss (same as maximizing the margin used by SVMs)

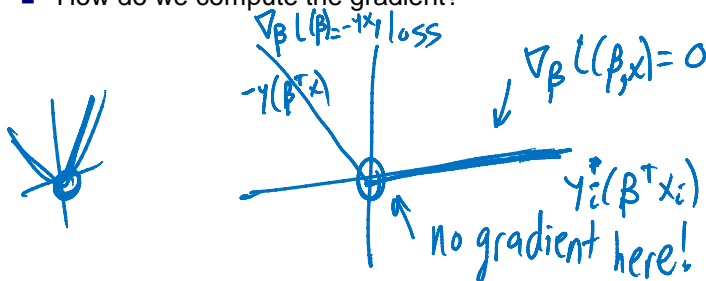


©Emily Fox 2014

27

## Minimizing Hinge Loss in Batch Setting

- Given a dataset:  $(x_1, y_1), \dots, (x_n, y_n)$
- Minimize average hinge loss:  $\min_{\beta} \frac{1}{n} \sum_{i=1}^n L(\beta, x_i)$
- How do we compute the gradient?

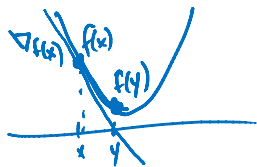


©Emily Fox 2014

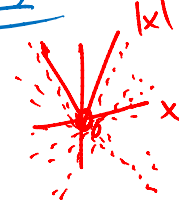
28

# Subgradients of Convex Functions

- Gradients lower bound convex functions:



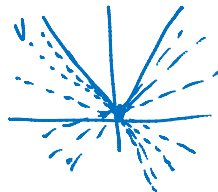
$$f(y) \geq f(x) + \nabla f(x)(y-x)$$



- Gradients are unique at  $x$  if function differentiable at  $x$

- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



For  $\|\beta_j\|:$   
 $\forall \beta \in [-1, 1]$

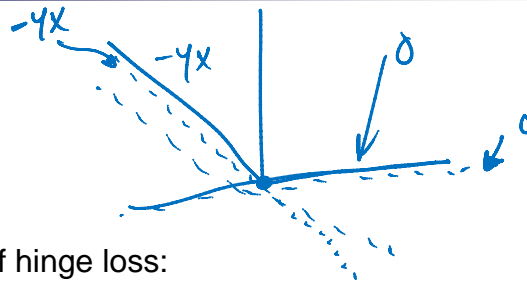
$\forall \beta \nabla f(x)$  subgradient  
 if  
 $f(y) \geq f(x) + \beta(y-x)$

©Emily Fox 2014

29

# Subgradient of Hinge

- Hinge loss:



- Subgradient of hinge loss:

- If  $y_t(\beta \cdot x_t) > 0$ :  $\nabla = 0$
- If  $y_t(\beta \cdot x_t) < 0$ :  $\nabla = -y_t x_t$
- If  $y_t(\beta \cdot x_t) = 0$ :  $\nabla = [-y_t x_t, 0]$
- In one line:

$$\nabla_{\beta} l(\beta, x) = \mathbb{I}(y(\beta^+ x) \leq 0) \cdot (-y x)$$

mistake indic.

©Emily Fox 2014

30

## Subgradient Descent for Hinge Minimization

- Given data:  $(x_1, y_1), \dots, (x_n, y_n)$

- Want to minimize:

$$\frac{1}{n} \sum_{i=1}^n \ell(\beta, x_i) = \frac{1}{n} \sum_{i=1}^n (-y_i(\beta \cdot x_i))_+$$

- Subgradient descent works the same as gradient descent: (reg. -  $\gamma_k$ )

- But if there are multiple subgradients at a point, just pick (any) one:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} - \underset{\substack{\uparrow \\ \text{step} \\ \text{size}}}{\eta} \sum_{i=1}^n \underset{\substack{\uparrow \\ \mathbb{I}(y_i(\beta^T x_i) \leq 0) \cdot (-y_i x_i)}}{\nabla \ell(\beta, x_i)}$$

©Emily Fox 2014

31

## Perceptron Revisited

- Perceptron update:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} + \underset{\substack{\uparrow \\ \text{if mistake}}}{\eta} \mathbb{I}[y_t(\beta^{(t)} \cdot x_t) \leq 0] \underset{\substack{\uparrow \\ \text{neg. subgrad.}}}{y_t x_t}$$

- Batch hinge minimization update:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} + \underset{\substack{\uparrow \\ \text{step.}}}{\eta} \frac{1}{n} \sum_{i=1}^n \left\{ \underset{\substack{\uparrow \\ \text{if mistake}}}{\mathbb{I}[y_i(\beta^{(t)} \cdot x_i) \leq 0]} \underset{\substack{\uparrow \\ \text{neg. subgrad.}}}{y_i x_i} \right\}$$

- Difference?

Perceptron algorithm = SGD  
for hinge loss minim. using  $\eta=1$ .

©Emily Fox 2014

32

# What you need to know

- Notion of online learning
- Perceptron algorithm
- Mistake bounds and proof (linearly separable)
- In online learning, report averaged weights at the end
- Perceptron is optimizing hinge loss
- Subgradients and hinge loss
- (Sub)gradient decent for hinge objective