

Module 5: Classification

Kernelized Perceptron: Kernels, again!

STAT/BIOSTAT 527, University of Washington

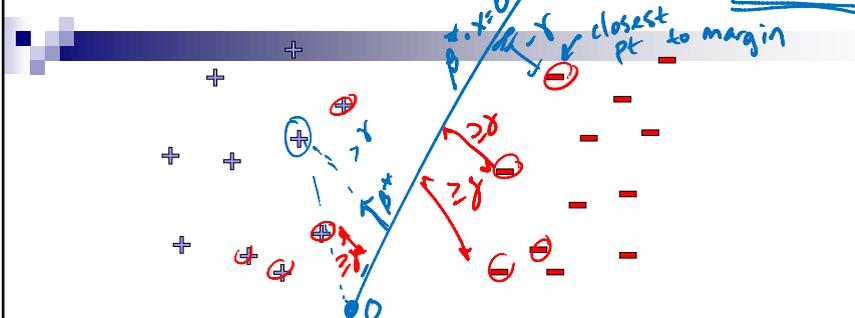
Emily Fox

May 29th, 2014

©Emily Fox 2014

1

Linear Separability: More formally, Using Margin



- Data linearly separable, if there exists

- a vector $\exists \beta^* \quad \|\beta^*\| = 1$
- a margin $\delta > 0$

- Such that

all pts are δ far away or more from $\beta^* \cdot x = 0$

$\forall t$ if $y_t = +1 \quad \beta^* \cdot x_t > \delta$
 $y_t = -1 \quad \beta^* \cdot x_t < -\delta \Rightarrow \boxed{y_t (\beta^* \cdot x_t) > \delta}$

©Emily Fox 2014

2

Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:

- Given a sequence of labeled examples: $(x_1, y_1), \dots, (x_n, y_n)$
examples need not be iid nor random
- Each covariate vector has bounded norm:

$$\forall t \quad \|x_t\| \leq R$$

- If dataset is linearly separable:

$$\exists \beta^* \quad (\beta^* \text{ is circled}) \quad \text{s.t.} \quad y_t (\beta^* \cdot x_t) \geq \gamma \quad \text{for some } \gamma > 0$$

- Then the number of mistakes made by the ~~online perceptron~~ on this sequence is bounded by

$$\left(\frac{R}{\gamma}\right)^2$$

crazy!

constant... doesn't depend on T

w.r.t. time

or dim X

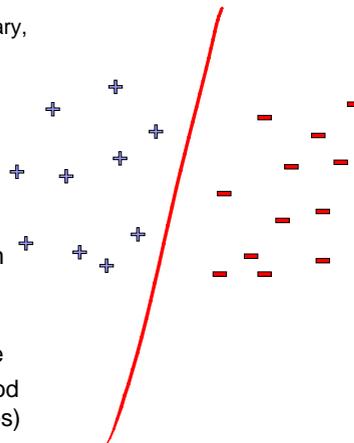
©Emily Fox 2014

3

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

4

What if the data are not linearly separable?

Use features of features of features of features....

$\phi(x) : \mathbb{R}^d \rightarrow F^d$

$x \in \mathbb{R}^d \rightarrow \phi(x) \in \mathbb{R}^{d'}$
 $d' \gg d$

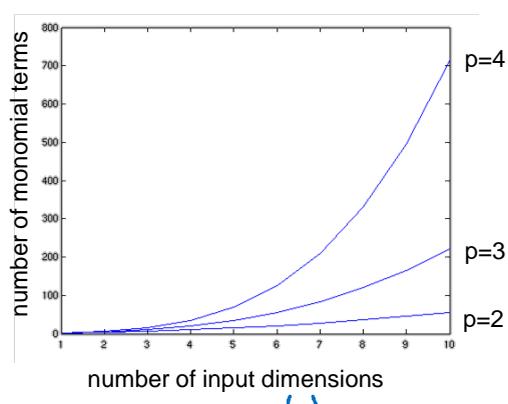
no plane separates data

Feature space can get really large really quickly!

Higher Order Polynomials *(basis expansion)*

$x \rightarrow \phi(x)$

$$\# \text{ terms} = \binom{p+d-1}{p} = \frac{(p+d-1)!}{p!(d-1)!}$$



d – covariate dimension
 p – degree of polynomial

even though dims of $\phi(x)$ are huge, we can fit the model very quickly.

grows fast!
 p = 6, d = 100
 about 1.6 billion terms

Perceptron Revisited

- Given weight vector $\beta^{(t)}$, predict point x by:

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

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

- Thus, write weight vector in terms of mistaken data points only:

- Let $M^{(t)}$ be time steps up to t when mistakes were made:

$$\beta^{(t)} = \sum_{j \in M^{(t)}} y_j x_j$$

- Prediction rule now:

$$\hat{y} = \text{sign}(\beta^{(t)} \cdot x) = \text{sign}\left(\sum_{j \in M^{(t)}} y_j x_j \cdot x\right) = \text{sign}\left(\sum_{j \in M^{(t)}} y_j (x_j \cdot x)\right)$$

- When using high dimensional features:

$$\hat{y} = \text{sign}\left(\sum_{j \in M^{(t)}} y_j (\phi(x_j) \cdot \phi(x))\right)$$

I compute efficiently?? classification only depends on inner product of x 's
when can

©Emily Fox 2014

7

Dot-Product of Polynomials

$\phi(u) \cdot \phi(v) =$ polynomials of degree exactly p

$$p=1: \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \cdot \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = u_1 v_1 + u_2 v_2 = (u \cdot v)^1$$

$$p=2: \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1^2 \\ u_1 u_2 \\ u_2 u_1 \\ u_2^2 \end{pmatrix} \cdot \begin{pmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{pmatrix} = u_1^2 v_1^2 + 2u_1 u_2 v_1 v_2 + u_2^2 v_2^2 = (u_1 v_1 + u_2 v_2)^2$$

By induction,

$$\text{degree } p: \phi(u) \cdot \phi(v) = (u \cdot v)^p$$

"kernel trick" compute efficiently!

©Emily Fox 2014

8

The Kernel Trick Again: Kernelized Perceptron

- Every time you make a mistake, remember (x_t, y_t)

keep track of indices $M^{(t)}$
where we make mistakes up to t .

- Kernelized perceptron prediction for x :

$$\begin{aligned} \text{sign}(\beta^{(t)} \cdot \phi(x)) &= \sum_{i \in M^{(t)}} y_i (\phi(x_i) \cdot \phi(x)) \\ &= \sum_{i \in M^{(t)}} y_i k(x_i, x) \end{aligned}$$

$(x_i \cdot x)$
 $k(x_i, x)$ ← kernel fcn.
 $k(u, v) = \phi(u) \cdot \phi(v)$
 one choice
 $(x_i \cdot x)^p$

©Emily Fox 2014

9

Going Infinite...

Change of notation:

$$h_j(x) \rightarrow \phi_j(x)$$

- Nonparametric Gaussian regression: $M \rightarrow \infty$

- Prior: $p(\beta | 0, \alpha^{-1} I_M)$

- Predictions: $f = \Phi \beta$

Dist. on f

$\begin{bmatrix} \phi(x_1) & \dots & \phi(x_n) \end{bmatrix}^T \beta$
 $\begin{bmatrix} \beta_1 & \dots & \beta_n \end{bmatrix}^T$
 linear comb. of Gaussians β_j Gauss.

$P(f) = N(f | 0, \alpha^{-1} \Phi \Phi^T)$
 $\Phi \Phi^T = \begin{bmatrix} \phi(x_1) \dots \phi(x_1) & \dots & \phi(x_1) \dots \phi(x_n) \\ \vdots & \ddots & \vdots \\ \phi(x_n) \dots \phi(x_1) & \dots & \phi(x_n) \dots \phi(x_n) \end{bmatrix}$

$E[f] = \Phi E[\beta] = 0$
 $\text{cov}(f) = \Phi E[\beta \beta^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T$

$\begin{bmatrix} \phi(x_1) & \dots & \phi(x_n) \end{bmatrix}$

What kernel can represent a $\phi \cdot \phi^T$??

- Gaussian process models replace explicit basis function representation with a direct specification in terms of a

positive definite kernel function

©Emily Fox 2014

10

Mercer Kernel Functions

- Predictions are of the form

$$p(f) = N(f | 0, \alpha^{-1} \Phi \Phi^T) = N(f | 0, K)$$

M = (1000, 10000, ...) *n x n matrix*

where the **Gram matrix** K is defined as

$$K_{ij} = k(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

kernel fcn *dim M*

- K is a **Mercer kernel** if the Gram matrix is positive definite for any n and any x_1, \dots, x_n

Note: K is $n \times n$ regardless of M (dim of basis/features)

Example of the "kernel trick"

©Emily Fox 2014

11

Polynomial kernels

- All monomials of degree d in $O(d)$ operations:

$$\phi(u) \cdot \phi(v) = (u \cdot v)^p = \text{polynomials of degree exactly } p$$

- How about all monomials of degree up to p ?

□ Solution 0: $\phi(u) \cdot \phi(v) = \sum_{i=0}^p \binom{p}{i} (u \cdot v)^i$

- Better solution:

$$(u \cdot v)^0 + (u \cdot v)^1 + (u \cdot v)^2 + \dots + (u \cdot v)^p = (u \cdot v + 1)^p$$

for degree p : $\phi(u) \cdot \phi(v) = k(u, v) = (u \cdot v + 1)^p$

©Emily Fox 2014

12

Common Kernels

- Polynomials of degree exactly d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$$

- Polynomials of degree up to d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$$

- Gaussian (squared exponential) kernel

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

(MANY MANY MORE!!)

radial
has a $\phi(\mathbf{u}) \cdot \phi(\mathbf{v})$
decomp.
 $\phi(\mathbf{u}) \in \mathbb{R}^\infty$
inner product
on an ∞ -dim.
space

What you need to know

- Notion of online learning
- Perceptron algorithm
- Mistake bounds and proofs
- The kernel trick
- Kernelized perceptron
- Derive polynomial kernel
- Common kernels
- In online learning, report averaged weights at the end

Module 5: Classification

Support Vector Machines

perceptron + regularization

STAT/BIOSTAT 527, University of Washington

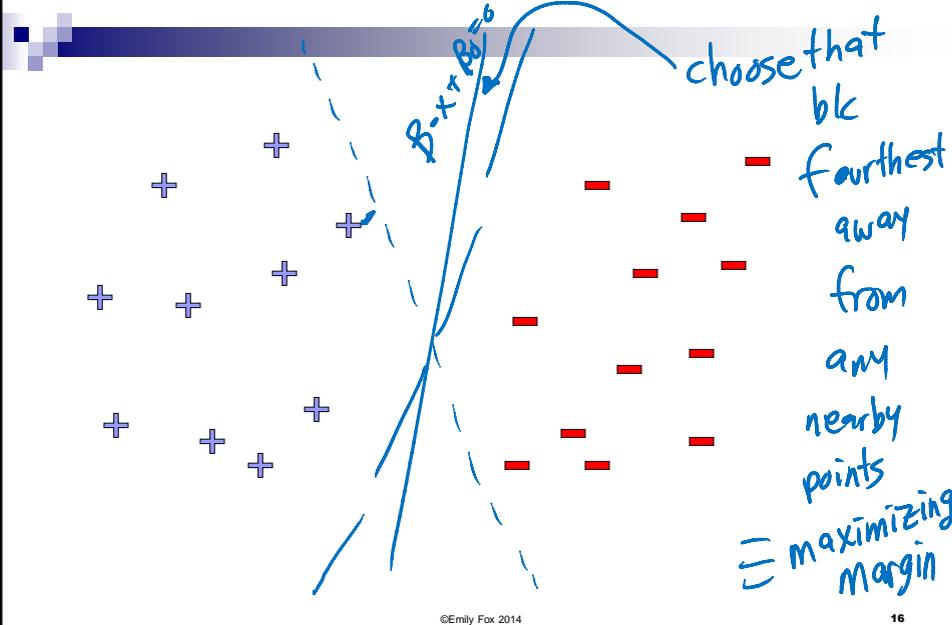
Emily Fox

May 29th, 2014

©Emily Fox 2014

15

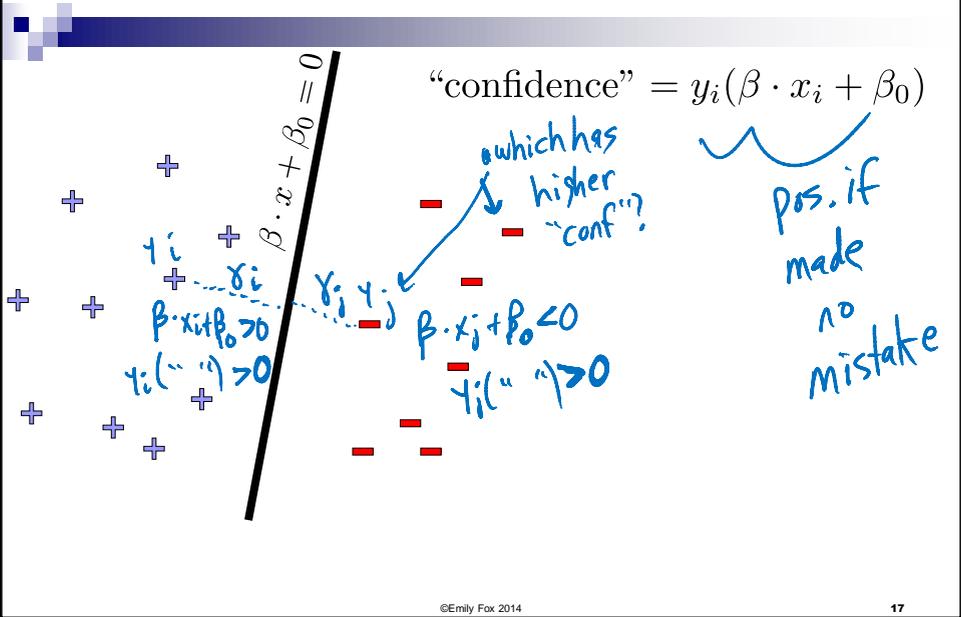
Linear classifiers – Which line is better?



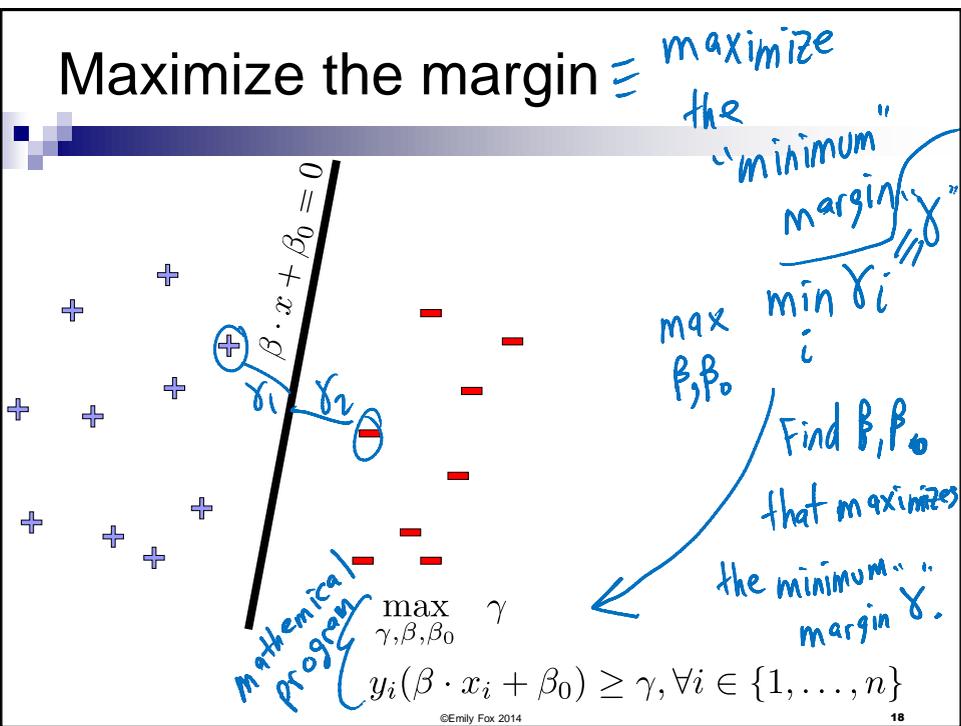
©Emily Fox 2014

16

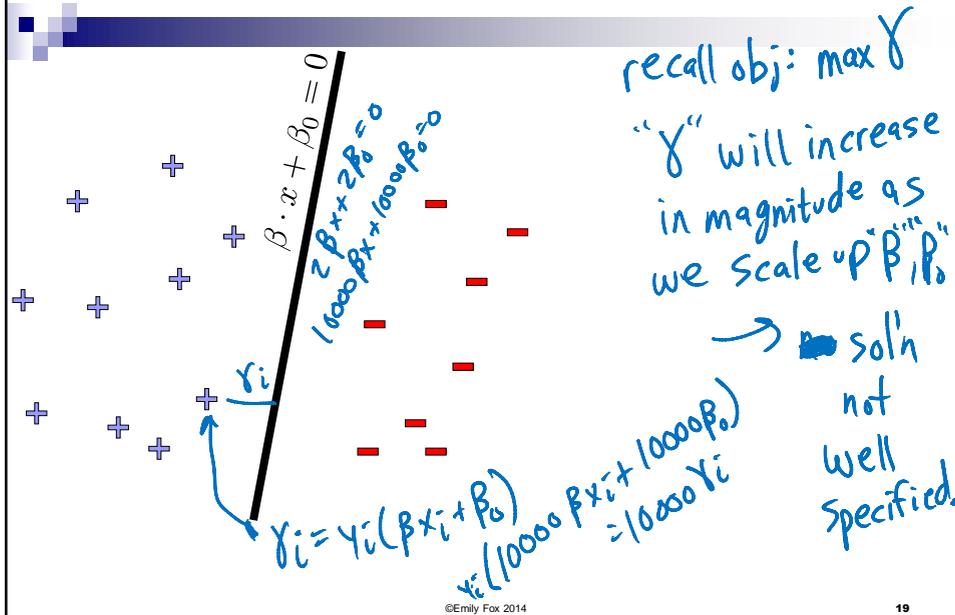
Pick the one with the largest margin!



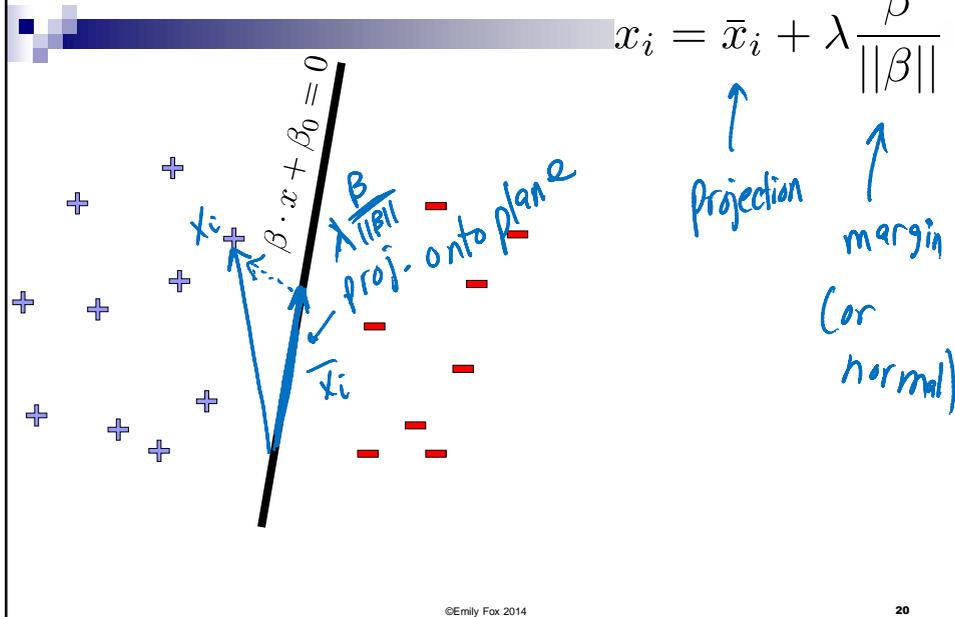
Maximize the margin \equiv maximize the “minimum margin”



But there are many planes...

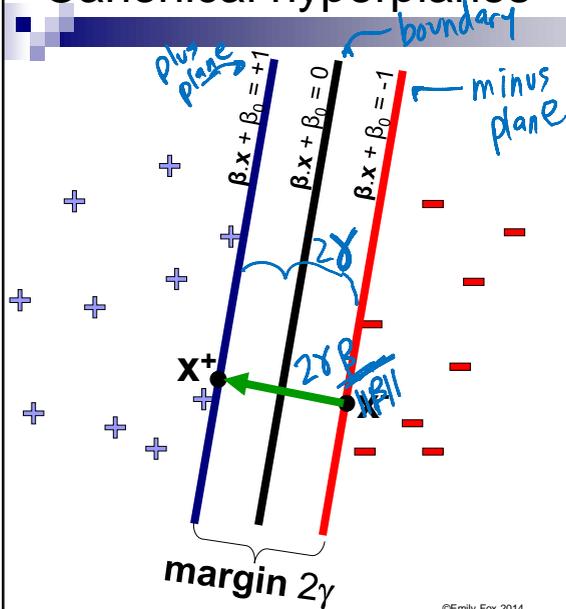


Review: Normal to a plane



A Convention: Normalized margin – Canonical hyperplanes

$$x_i = \bar{x}_i + \lambda \frac{\beta}{\|\beta\|}$$



pts. on plane

$$\beta \cdot x^+ + \beta_0 = +1$$

$$\beta \cdot x^- + \beta_0 = -1$$

$$x^+ = x^- + 2\gamma \frac{\beta}{\|\beta\|}$$

$$\beta \cdot (x^- + 2\gamma \frac{\beta}{\|\beta\|}) + \beta_0 = +1$$

$$\beta \cdot x^- + \beta_0 + 2\gamma \frac{\beta \cdot \beta}{\|\beta\|} = +1$$

$$\underbrace{-1}_{-1} + 2\gamma \frac{1}{\|\beta\|} = +1$$

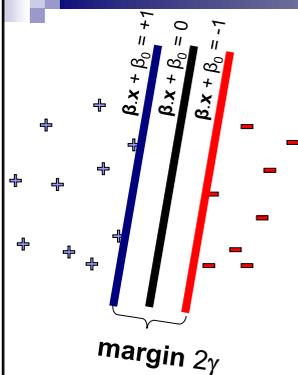
$$\rightarrow 2\gamma = \frac{1}{\|\beta\|}$$

©Emily Fox 2014

21

Margin maximization using canonical hyperplanes

$$\max_x \frac{1}{f(x)} \Leftrightarrow \min_x f(x)$$



Unnormalized problem: $\max_{\gamma, \beta, \beta_0} \gamma$ "min" margin

$$y_i(\beta \cdot x_i + \beta_0) \geq \gamma, \forall i \in \{1, \dots, n\}$$

Normalized Problem:

if $y_i(\beta \cdot x_i + \beta_0) \geq 1$

$\max \gamma \Leftrightarrow$

$\max \frac{1}{\|\beta\|}$

"canonical" hyperplanes

$\Leftrightarrow \min \|\beta\|^2$

standard SVM (derivatives problem are easier) (norm minimization)

$$\min_{\beta, \beta_0} \|\beta\|_2^2$$

$$y_i(\beta \cdot x_i + \beta_0) \geq 1, \forall i \in \{1, \dots, n\}$$

©Emily Fox 2014

22

Support vector machines (SVMs)

$$\min_{\beta, \beta_0} \|\beta\|_2^2$$

$$y_i(\beta \cdot x_i + \beta_0) \geq 1, \forall i \in \{1, \dots, n\}$$

- Solve efficiently by many methods, e.g.,
 - quadratic programming (QP)
 - Well-studied solution algorithms
 - Stochastic gradient descent
- Hyperplane defined by support vectors

** only perturbations of these points will change the sol'n. (contrast w/ LMA, QDA)*

©Emily Fox 2014 23

What if the data are not linearly separable?

Use features of features of features of features....

SVM can be solved as usual BUT now we can apply a kernel (just as in perceptron)

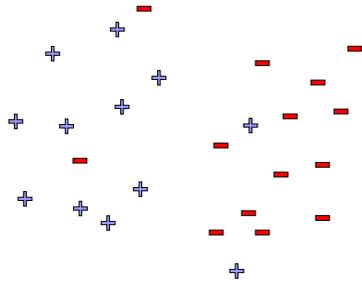
©Emily Fox 2014 24

What if the data are still not linearly separable?



$$\min_{\beta, \beta_0} \|\beta\|_2^2$$

$$y_i(\beta \cdot x_i + \beta_0) \geq 1, \forall i \in \{1, \dots, n\}$$



- If data are not linearly separable, some points don't satisfy margin constraint:
- How bad is the violation?
- Tradeoff margin violation with $\|\beta\|$:

©Emily Fox 2014

25

SVMs for Non-Linearly Separable meet my friend the Perceptron...



- Perceptron was minimizing the hinge loss:

$$\sum_{i=1}^n (-y_i(\beta \cdot x_i + \beta_0))_+$$

- SVMs minimizes the regularized hinge loss!!

$$\|\beta\|_2^2 + C \sum_{i=1}^n (1 - y_i(\beta \cdot x_i + \beta_0))_+$$

©Emily Fox 2014

26

Stochastic Gradient Descent for SVMs

- Perceptron minimization:

$$\sum_{i=1}^n (-y_i(\beta \cdot x_i + \beta_0))_+$$

- SGD for Perceptron:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} + \mathbb{I}[y_t(\beta^{(t)} \cdot x_t) \leq 0] y_t x_t$$

- SVMs minimization:

$$\|\beta\|_2^2 + C \sum_{i=1}^n (1 - y_i(\beta \cdot x_i + \beta_0))_+$$

- SGD for SVMs:

©Emily Fox 2014

27

Side note: What's the difference between SVMs and logistic regression?

SVM:

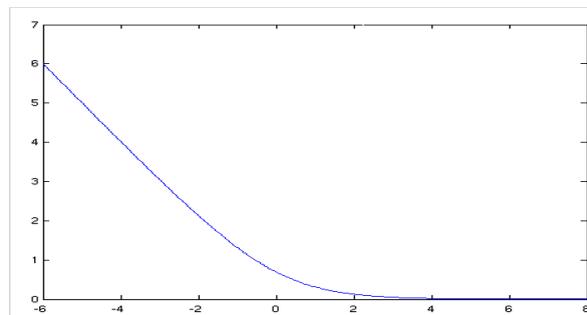
$$\min_{\beta, \beta_0} \|\beta\|_2^2 + C \sum_{i=1}^n (1 - y_i(\beta \cdot x_i + \beta_0))_+$$

Logistic regression:

$$p(Y = 1 | x, \beta) = \frac{1}{1 + e^{-(\beta \cdot x + \beta_0)}}$$

Log loss:

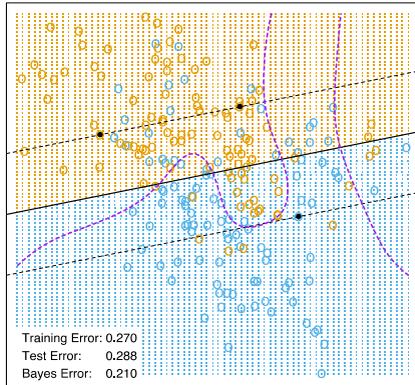
$$-\log p(Y = 1 | x, \beta) = \log(1 + e^{-(\beta \cdot x + \beta_0)})$$



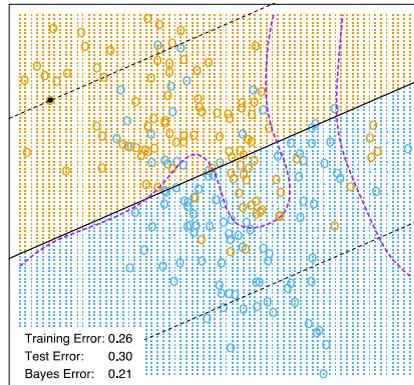
©Emily Fox 2014

28

Mixture Model Example



$C = 10000$



$C = 0.01$

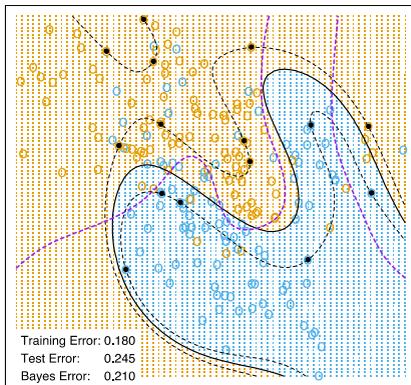
From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

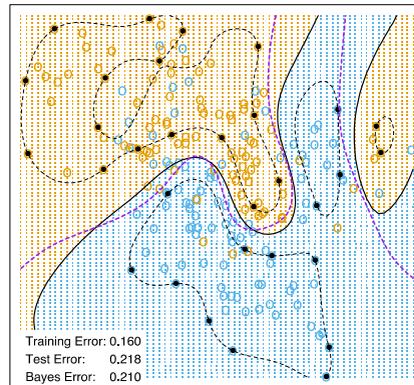
29

Mixture Model Example - Kernels

SVM - Degree-4 Polynomial in Feature Space



SVM - Radial Kernel in Feature Space



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

30