

Module 5: Classification

Kernelized Perceptron: Kernels, again!

STAT/BIOSTAT 527, University of Washington

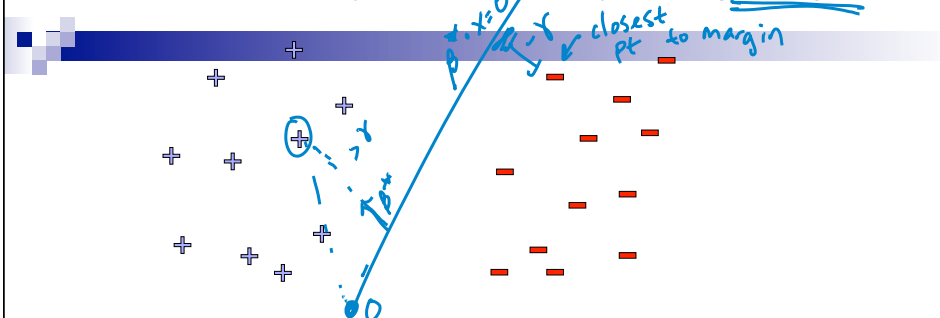
Emily Fox

June 4th, 2013

©Emily Fox 2013

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 $\beta^* \cdot x = 0$

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

©Emily Fox 2013

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:

$$\exists \beta^* \quad \|\beta^*\| = 1 \quad \text{s.t.} \quad y_t (\beta^* \cdot x_t) \geq \gamma \quad \text{for some } \gamma > 0$$

- 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 or dim X

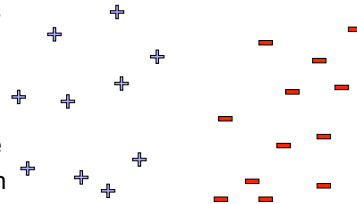
©Emily Fox 2013

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



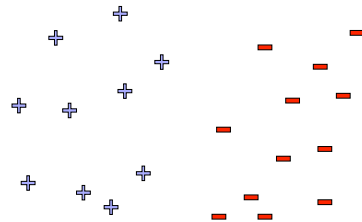
- However, real world not linearly separable

- 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 2013

4

What if the data are not linearly separable?



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

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

Feature space can get really large really quickly!

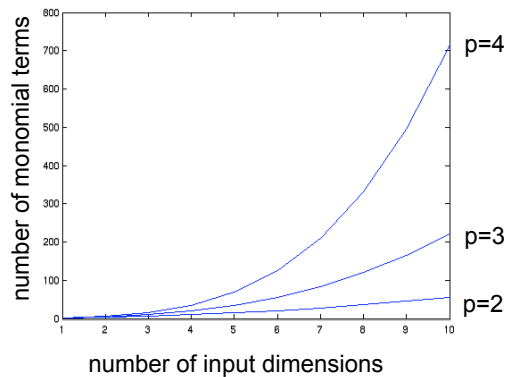
©Emily Fox 2013

Higher Order Polynomials



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

d – covariate dimension
p – degree of polynomial



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

©Emily Fox 2013

6

Perceptron Revisited

- Given weight vector $\beta^{(t)}$, predict point \mathbf{x} by:
- Mistake at time t : $\beta^{(t+1)} \leftarrow \beta^{(t)} + y_t \mathbf{x}_t$
- Thus, write weight vector in terms of mistaken data points only:
 - Let $M^{(t)}$ be time steps up to t when mistakes were made:
- Prediction rule now:
- When using high dimensional features:

©Emily Fox 2013

7

Dot-Product of Polynomials

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

©Emily Fox 2013

8

The Kernel Trick Again: Kernelized Perceptron

- Every time you make a mistake, remember (x_t, y_t)
- Kernelized perceptron prediction for \mathbf{x} :

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

©Emily Fox 2013

9

★ Going Infinite...

Change of notation:

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

- Nonparametric Gaussian regression: *basis fcn's*
Would like to let the number of "features" $M \rightarrow \infty$

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

- Predictions: $f = \Phi \beta$

Dist. on f

linear comb. of Gaussians β_j

$\Rightarrow f$ Gauss.

$E[f] = \Phi E[\beta] = 0$

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

$P(f) = N(f \mid 0, \alpha^{-1} \Phi \Phi^T)$

$\begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$

$(\phi_1(x_1) \dots \phi_m(x_2))$

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

©Emily Fox 2013

10

Mercer Kernel Functions

- Predictions are of the form

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

$\alpha^{-1} \Phi \Phi^T$ is an $n \times n$ matrix
 Φ is $n \times M$, Φ^T is $M \times n$

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

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

$k(x_i, x_j)$ is the kernel function, $\phi(x_i)$ and $\phi(x_j)$ are vectors of dimension 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"

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:

□ Better solution:

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

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

- Sigmoid

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

©Emily Fox 2013

13

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

©Emily Fox 2013

14

Module 5: Classification

Support Vector Machines

STAT/BIOSTAT 527, University of Washington

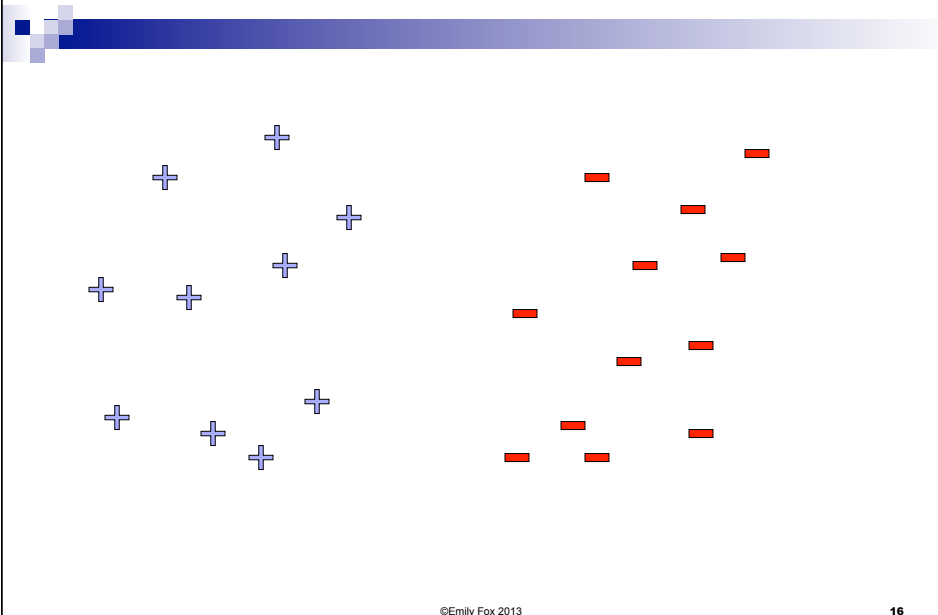
Emily Fox

June 4th, 2013

©Emily Fox 2013

15

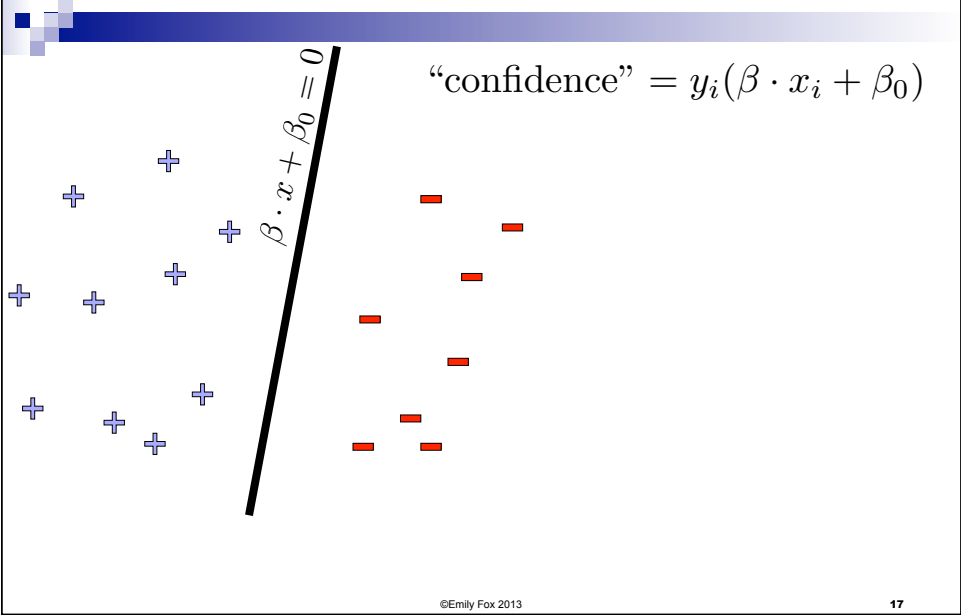
Linear classifiers – Which line is better?



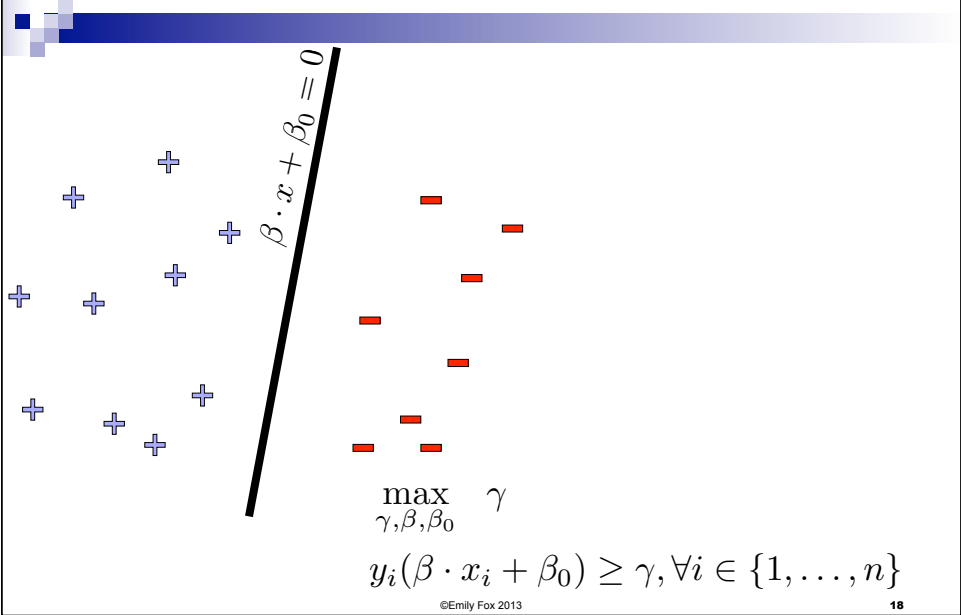
©Emily Fox 2013

16

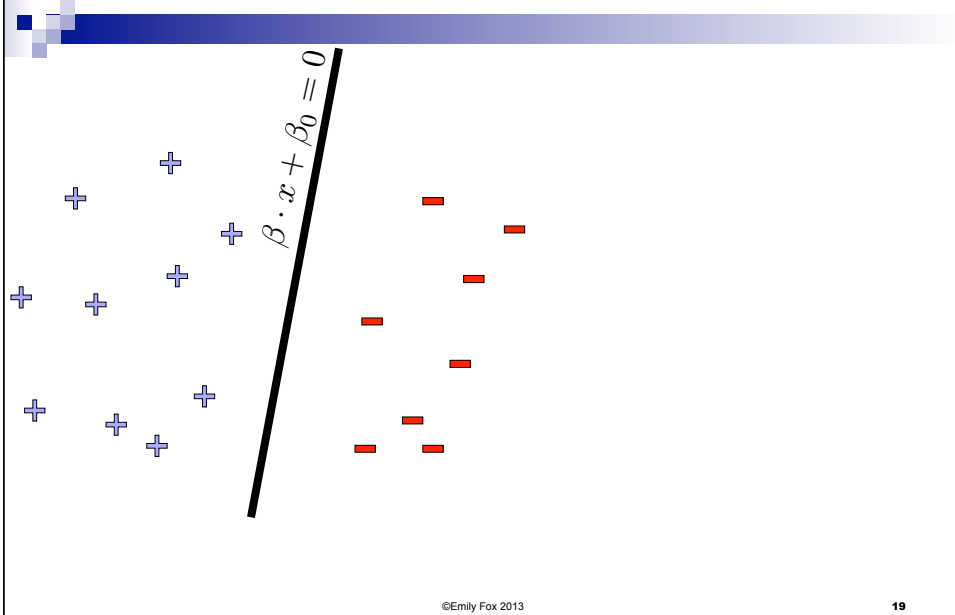
Pick the one with the largest margin!



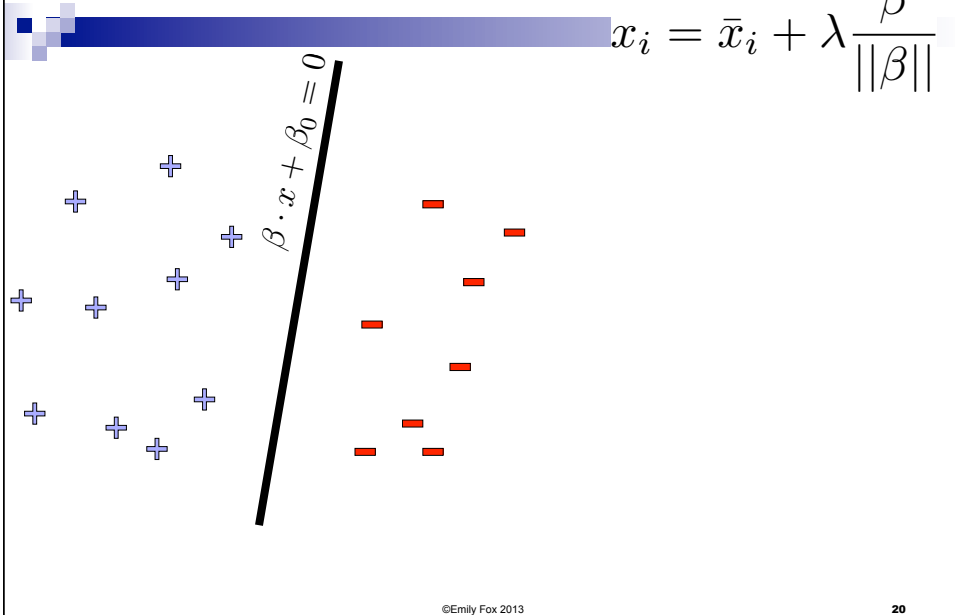
Maximize the margin



But there are many planes...

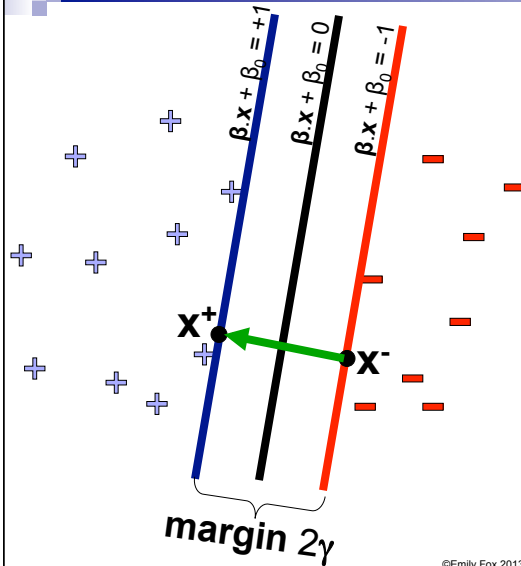


Review: Normal to a plane



A Convention: Normalized margin – Canonical hyperplanes

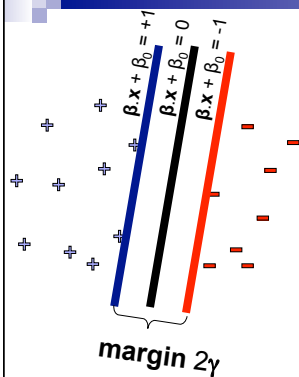
$$x_i = \bar{x}_i + \lambda \frac{\beta}{\|\beta\|}$$



©Emily Fox 2013

21

Margin maximization using canonical hyperplanes



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

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

Normalized Problem:

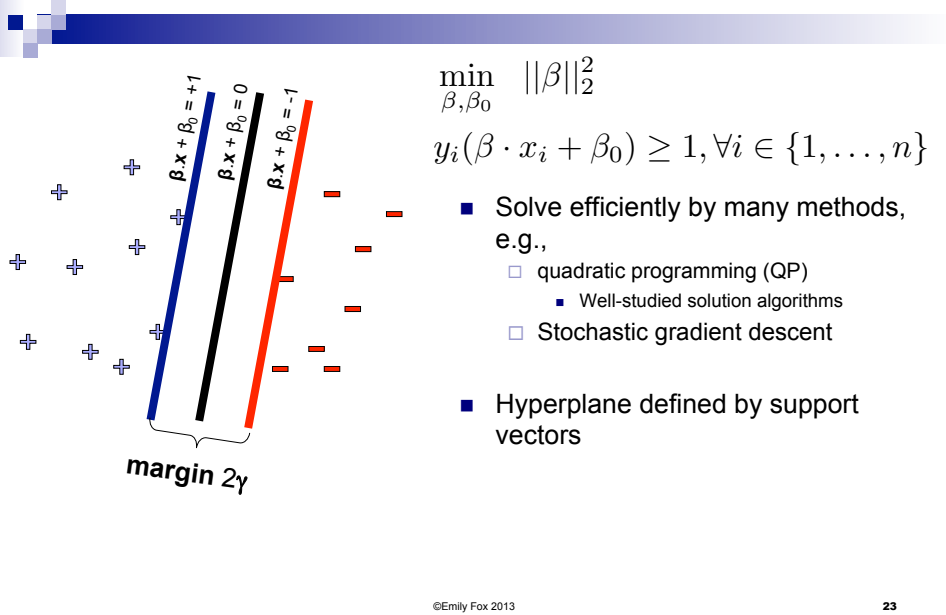
$$\min_{\beta, \beta_0} \|\beta\|_2^2$$

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

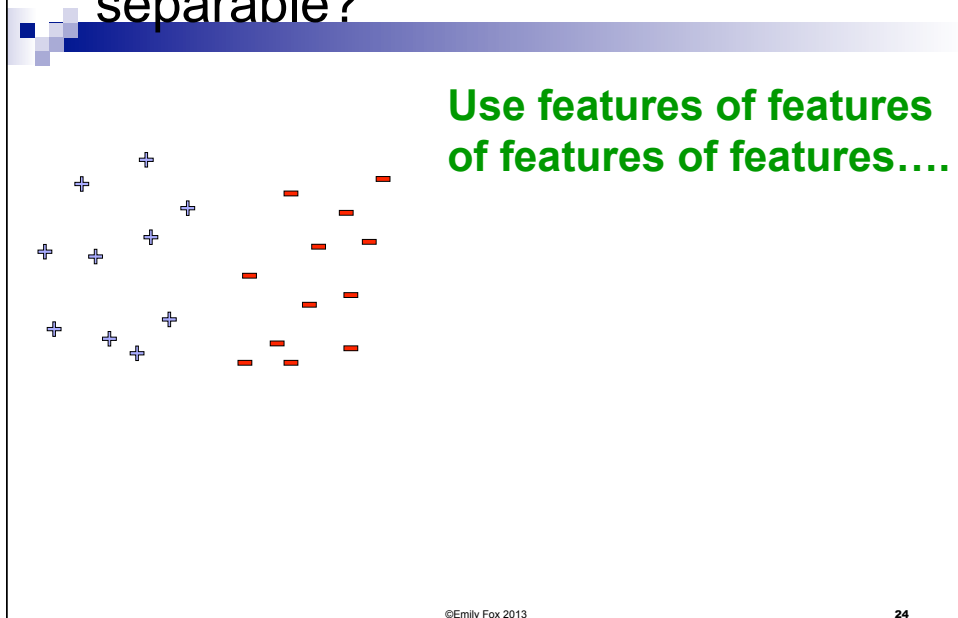
©Emily Fox 2013

22

Support vector machines (SVMs)



What if the data are not linearly separable?

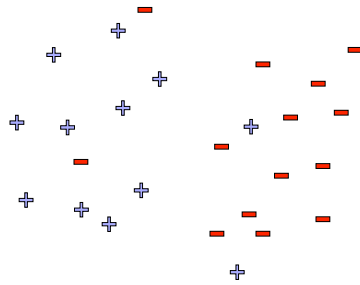


What if the data are still not linearly separable?



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

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



- 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 2013

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 2013

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 2013

27

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

SVM:

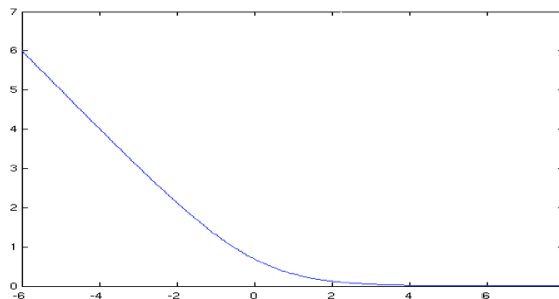
$$\min_{\beta, \beta_0} \|\beta\|_2^2 + C \sum_{i=1}^n (1 - y_i(\beta \cdot x_i + \beta_0))_+$$

Logistic regression:

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

Log loss:

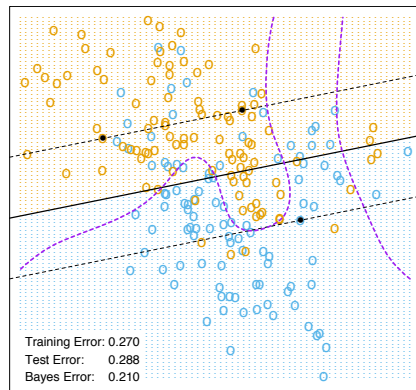
$$-\log p(Y = 1 | x, \beta) = \log(1 + e^{-(\beta \cdot x + \beta_0)})$$



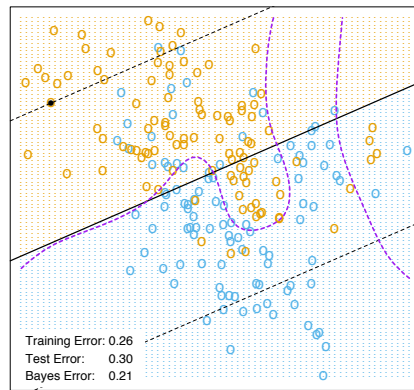
©Emily Fox 2013

28

Mixture Model Example



$C = 10000$



$C = 0.01$

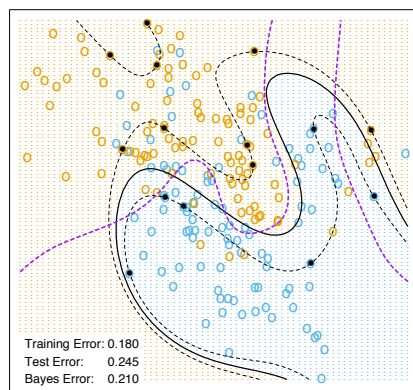
From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

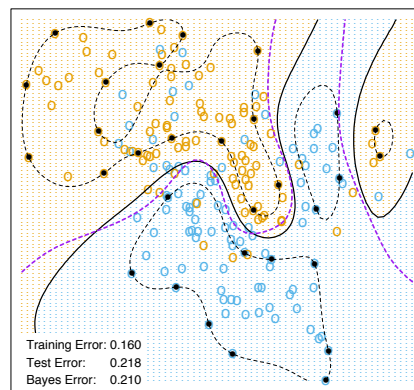
29

Mixture Model Example - Kernels

SVM - Degree-4 Polynomial in Feature Space



SVM - Radial Kernel in Feature Space

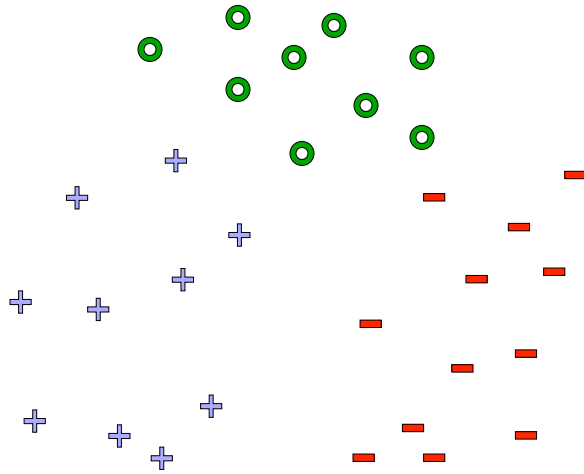


From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

30

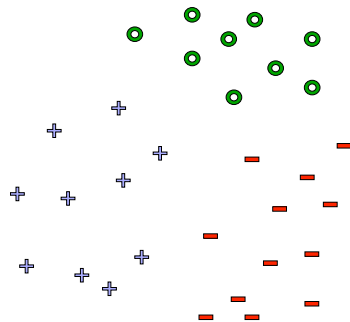
What about multiple classes?



©Emily Fox 2013

31

One against All



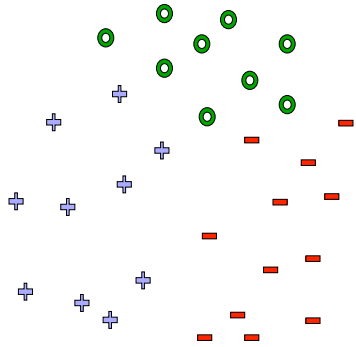
Learn 3 classifiers:

©Emily Fox 2013

32

Learn 1 classifier: Multiclass SVM

Simultaneously learn 3 sets of weights



$$\beta^{(y_i)} \cdot x_i + \beta_0^{(y_i)} \geq \beta^{(y')} \cdot x_i + \beta_0^{(y')} + 1, \forall y' \neq y_i, \forall i$$

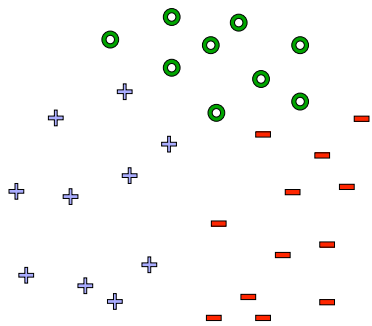
©Emily Fox 2013

33

Learn 1 classifier: Multiclass SVM

$$\beta^{(y_i)} \cdot x_i + \beta_0^{(y_i)} \geq \beta^{(y')} \cdot x_i + \beta_0^{(y')} + 1, \forall y' \neq y_i, \forall i$$

$$\min_{\beta, \beta_0} \sum_y \beta^{(y)} \cdot \beta^{(y)} + C \sum_i \xi_i$$



©Emily Fox 2013

34

What you need to know

- Maximizing margin
- Derivation of SVM formulation
- Non-linearly separable case
 - Hinge loss
 - A.K.A. adding slack variables
- SVMs = Perceptron + L2 regularization
- Can optimize SVMs with SGD
 - Many other approaches possible
- Handling multiple classes

Reading

- Hastie, Tibshirani, Friedman – 12.1-12.3