

Module 5: Classification

Kernelized Perceptron: Kernels, again!

STAT/BIOSTAT 527, University of Washington

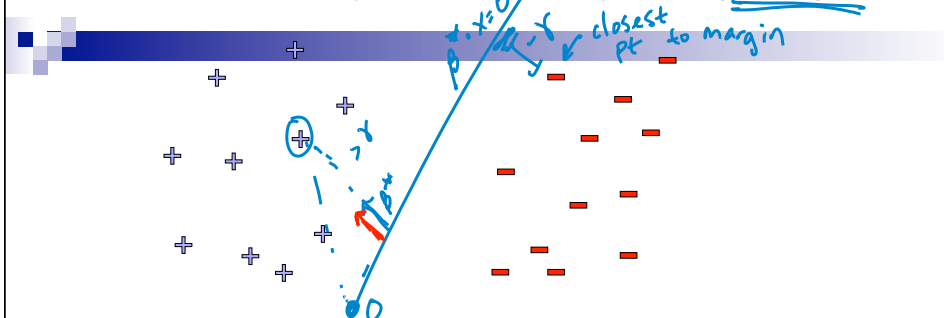
Emily Fox

June 4th, 2013

©Emily Fox 2013

1

Linear Separability: More formally, Using Margin



- Data linearly separable, if there exists

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

- 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

*If mistake
else $\beta^{(t+1)} \leftarrow \beta^{(t)} + y_t x_t$
 $\beta^{(t+1)} \leftarrow \beta^{(t)}$*

← t

for some $\gamma > 0$

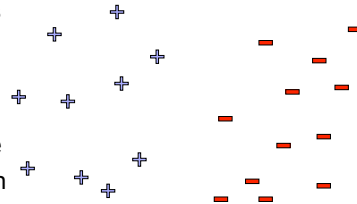
©Emily Fox 2013

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



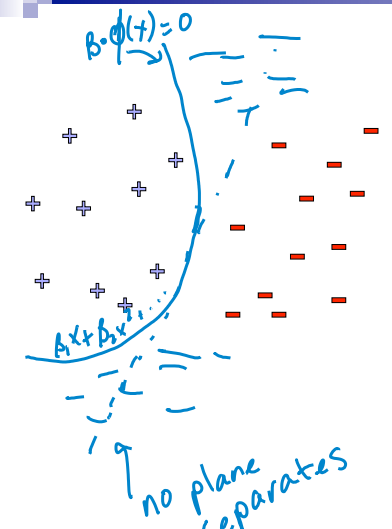
- However, real world not linearly separable

- 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....

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

could even be ∞ -dim

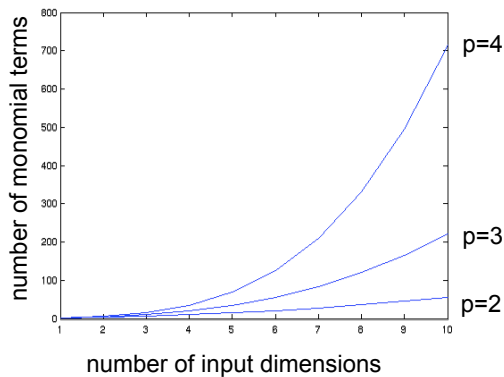
$$\phi(x) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ x_1 x_2 \\ e^{\sin x} \\ e^{-x_2^2} \\ \vdots \end{pmatrix}, \quad \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ \vdots \end{pmatrix}$$

space here

Feature space can get really large really quickly!

Higher Order Polynomials *as basis expansion*

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



d – covariate dimension
p – degree of polynomial

even though dims of $\phi(x)$ are huge, can fit model very, 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 : $\beta^{(t+1)} \leftarrow \beta^{(t)} + \underline{y_t x_t}$
- Thus, write weight vector in terms of mistaken data points only:
 - Let $M^{(t)}$ be time steps up to t when mistakes were made:

$$\beta^{(t)} = \sum_{j \in M^{(t)}} y_j x_j \quad (\text{assuming } \beta^{(0)} = 0)$$

- Prediction rule now:

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

- When using high dimensional features:

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

when can I compute this efficiently?
classification only depends on inner prod of x prod

©Emily Fox 2013

7

Dot-Product of Polynomials

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

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

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

pf by induction:

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

can compute very efficiently
 "kernel trick"

©Emily Fox 2013

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 :

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

*Kernel fcn
 $k(u, v) = \phi(u) \cdot \phi(v)$*

©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: $p(\beta \mid 0, \alpha^{-1} I_M)$ *dim M*

- Predictions: $f = \Phi \beta$

Dist. on f
 $\begin{pmatrix} f(x_1) & \dots & f(x_n) \end{pmatrix}^T = \begin{pmatrix} \phi_1(x_1) & \dots & \phi_M(x_1) \\ \vdots & & \vdots \\ \phi_1(x_n) & \dots & \phi_M(x_n) \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_M \end{pmatrix}$
linear comb. of Gaussians β_j
 $\Rightarrow f$ Gauss.
 $E[f] = \Phi E[\beta] = 0$
 $\text{cov}(f) = \Phi E[\beta \beta^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T$
 $\Phi = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_M(x_1) \\ \vdots & & \vdots \\ \phi_1(x_n) & \dots & \phi_M(x_n) \end{bmatrix}$
 $\begin{matrix} n \times 1 \\ \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \\ n \times M \end{matrix} \begin{matrix} M \times 1 \\ \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \\ M \times 1 \end{matrix}$

- 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 fcn, $\phi(x_i)^T \phi(x_j)$ is $\dim M$

- 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"

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: $d=2$
 $(u \cdot v)^1 + (u \cdot v)^2 + (v \cdot u)^1 + (u \cdot v)^0 = (u \cdot v + 1)^2$

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

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

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

- Strings
- Graph

eq. to $\phi(\mathbf{u}) \cdot \phi(\mathbf{v})$
dim of $\phi(\mathbf{u})$
 ∞
inner prod in
an ∞ -dim space
MANY, MANY
MORE...

©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

≈ perceptron + regularization

STAT/BIOSTAT 527, University of Washington

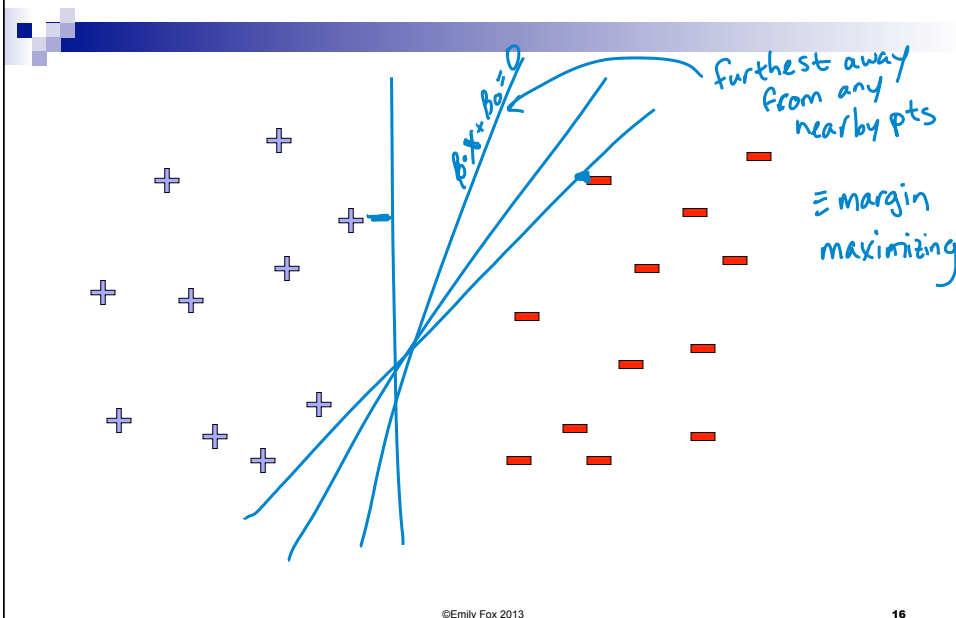
Emily Fox

June 4th, 2013

©Emily Fox 2013

15

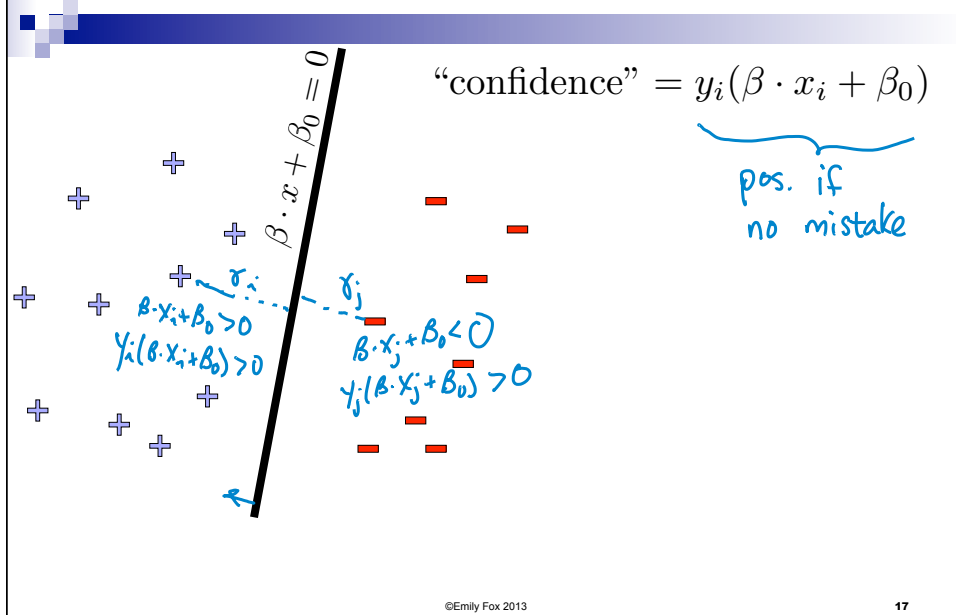
Linear classifiers – Which line is better?



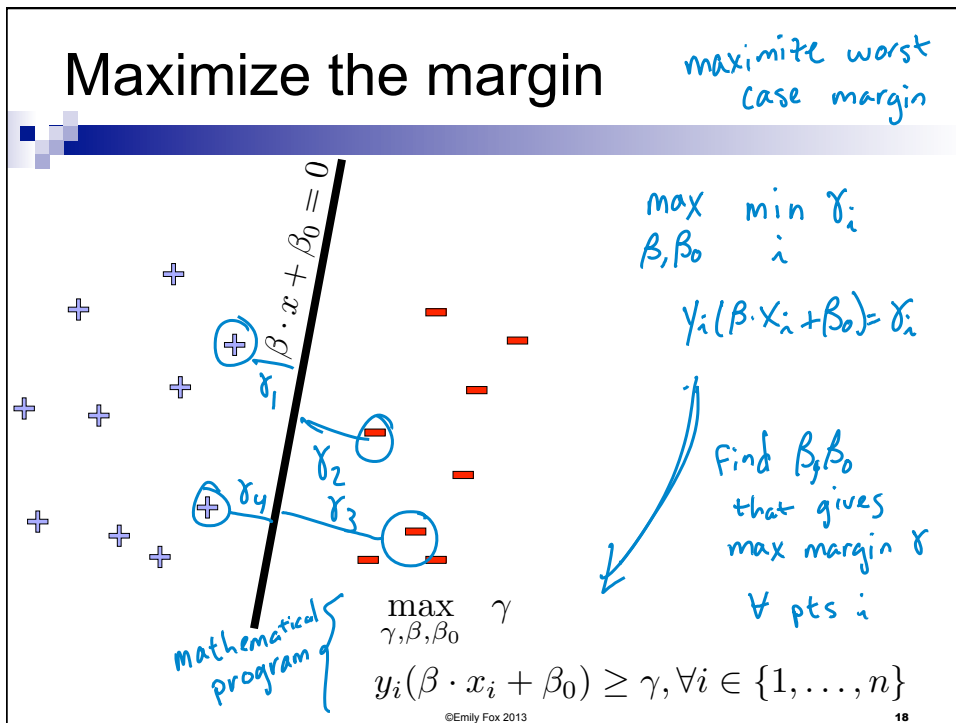
©Emily Fox 2013

16

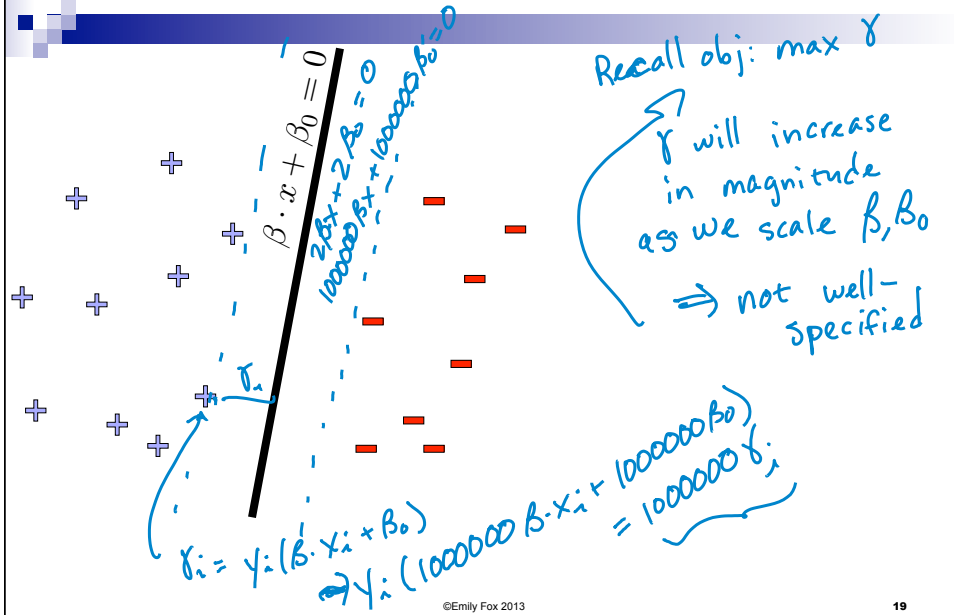
Pick the one with the largest margin!



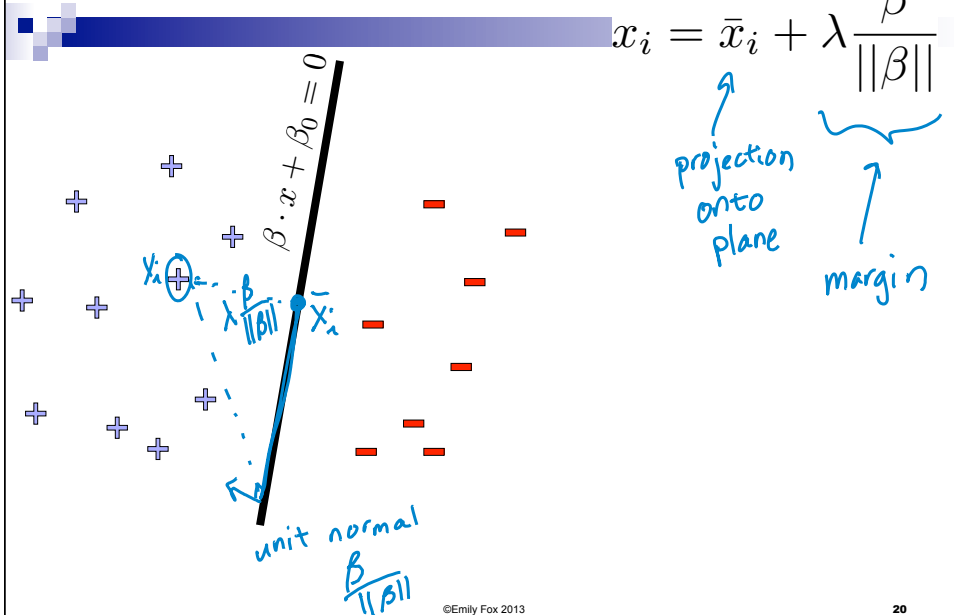
Maximize the margin



But there are many planes...

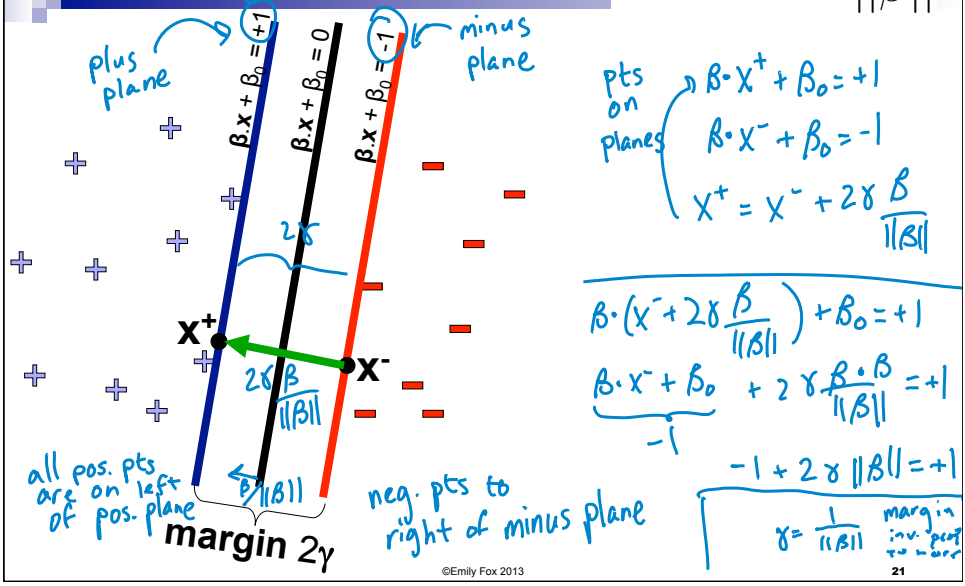


Review: Normal to a plane



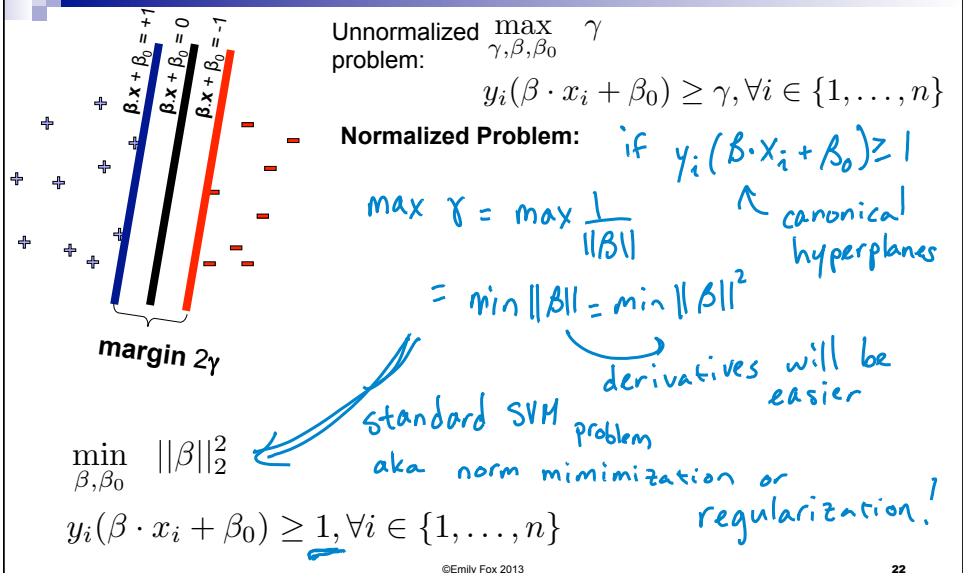
A Convention: Normalized margin – Canonical hyperplanes

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



Margin maximization using canonical hyperplanes

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



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 pts will change the soln (contrast w/ LDA, QDA)

©Emily Fox 2013 23

What if the data are not linearly separable?

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

$$\phi(x) = \begin{pmatrix} x^2 \\ x \\ e^x \sin x \\ \vdots \end{pmatrix}$$

can be done efficiently w/ kernels (just like in perceptron case)

SVMs popularized kernels in ML

too much var.?

(run risk of over-fitting)

©Emily Fox 2013 24

What if the data are still not linearly separable?

Now support vectors are
- on margin
or
- misclassified pts

$\min_{\beta, \beta_0} \|\beta\|_2^2$
 $y_i(\beta \cdot x_i + \beta_0) \geq 1, \forall i \in \{1, \dots, n\}$

- If data are not linearly separable, some points don't satisfy margin constraint:
 - $\exists i$ s.t. $y_i(\beta \cdot x_i + \beta_0) < 1$
 - or $1 - y_i(\beta \cdot x_i + \beta_0) > 0$ like hinge loss
- How bad is the violation?
 - constraint violation = $\begin{cases} 0 & \text{if } 1 - y_i(\beta \cdot x_i + \beta_0) \leq 0 \\ 1 - y_i(\beta \cdot x_i + \beta_0) & \text{otherwise} \end{cases}$
- Tradeoff margin violation with $\|\beta\|$:
 - $\min_{\beta, \beta_0} \|\beta\|_2^2 + C \sum_{i=1}^n (1 - y_i(\beta \cdot x_i + \beta_0))_+$
 - SVM objective

violates (with arrow pointing to the constraint equation)

©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!!

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

reg. (arrow to $\|\beta\|_2^2$)
reg. of fit (arrow to C)
just convention (arrow to \sum)

©Emily Fox 2013 26

$C=0 \rightarrow \beta=0 \rightarrow \text{ignore data}$

$C=\infty \rightarrow \text{want to separate data}$

Stochastic Gradient Descent for SVMs

- Perceptron minimization:

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

- SGD for Perceptron:

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

↑
update weights

make mistake

stepsize $\eta=1$

- kernelize like perceptron

- SVMs minimization:

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

- SGD for SVMs:

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

$$+ \eta (C \mathbb{I}[(1 - y_t(\beta \cdot x_t + \beta_0)) > 0] \cdot y_t x_t - 2\beta^{(t)})$$

need to play w/ stepsize

- pick C, η by CV

other ways, too

©Emily Fox 2013

27

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

SVM:

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

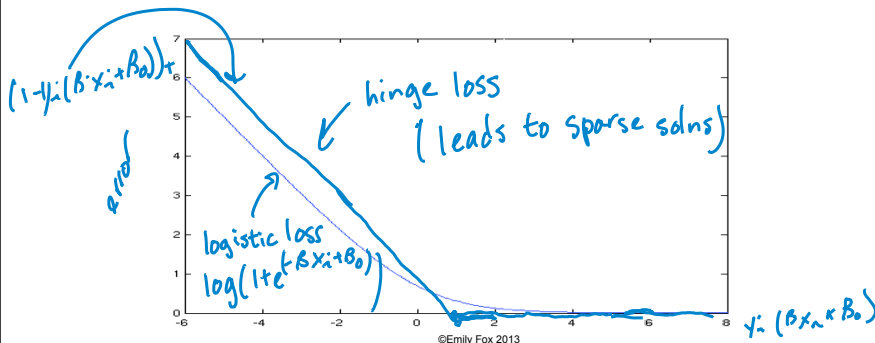
← only diff is loss fun

Logistic regression:

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

Log loss:

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

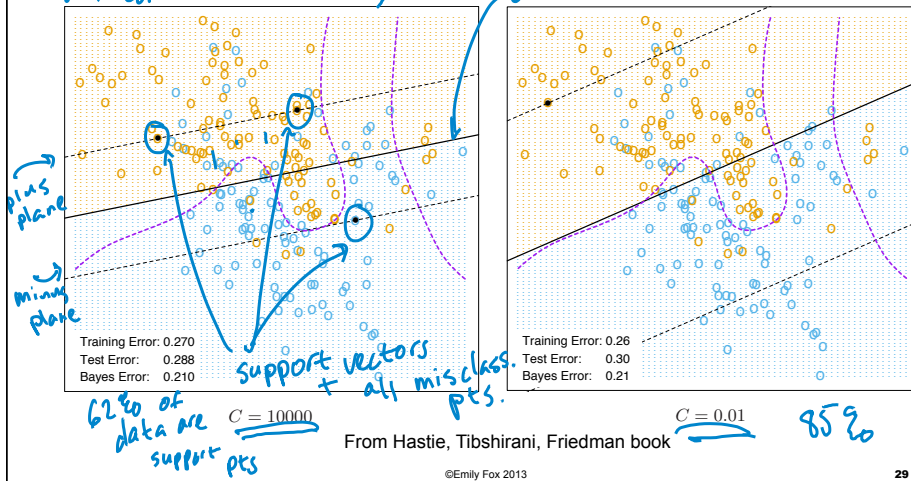


©Emily Fox 2013

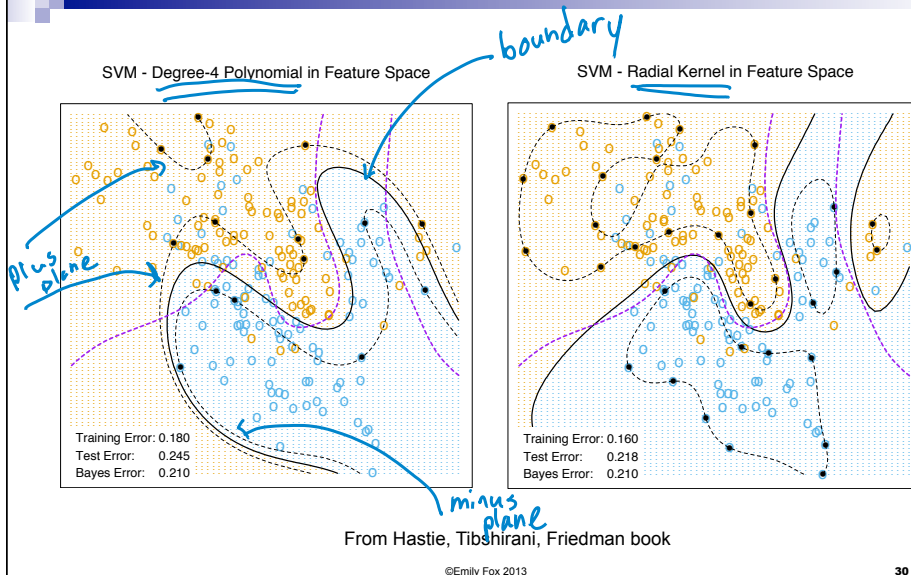
28

Mixture Model Example

2-class mixture model w/ overlap
still constrained to line, so not too affected by choice of C
decision boundary



Mixture Model Example - Kernels



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