

# Kernels

Machine Learning – CSE446

Carlos Guestrin

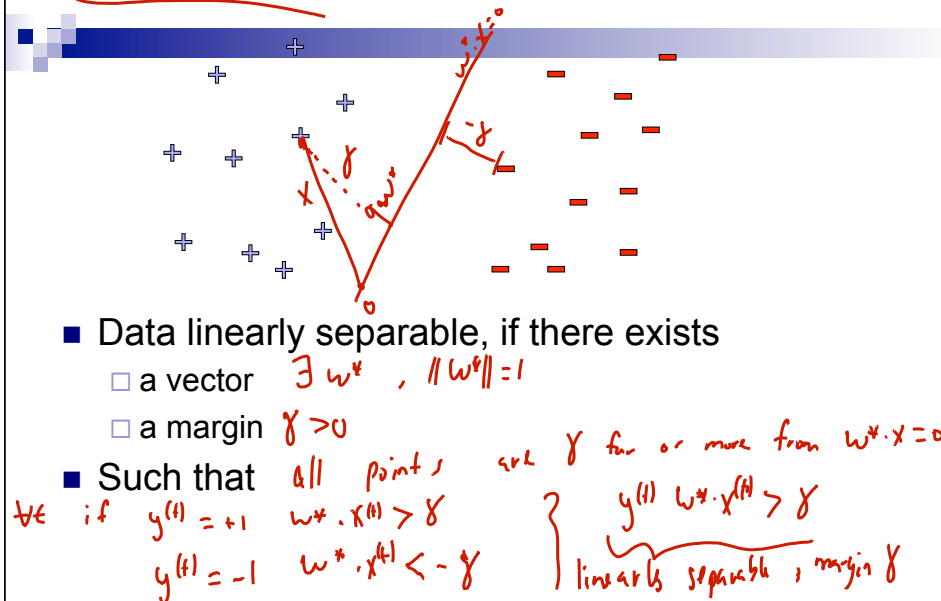
University of Washington

May 3, 2013

©Carlos Guestrin 2005-2013

1

## Linear Separability: More formally, Using Margin



©Carlos Guestrin 2005-2013

2

## Perceptron Analysis: Linearly Separable Case

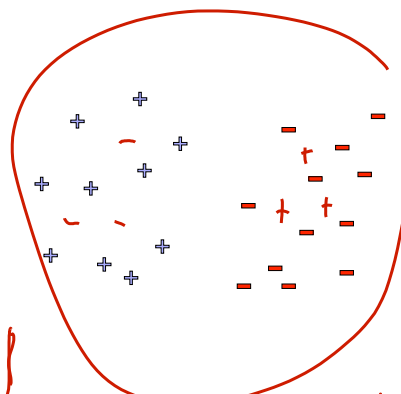
- Theorem [Block, Novikoff]:
  - Given a sequence of labeled examples:  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(T)}, y^{(T)})$   
*examples need not be iid. or random...*
  - Each feature vector has bounded norm:  
 $\forall t \quad \|x^{(t)}\| \leq R$   *$w^*$  is unknown!*
  - If dataset is linearly separable:  
 $\exists w^*, \|w^*\|=1 \quad \forall t \quad y^{(t)} w^* \cdot x^{(t)} \geq \gamma$ , for  $\gamma > 0$
- Then the number of mistakes made by the online perceptron on ~~any~~ this sequence is bounded by  
 $\left(\frac{R}{\gamma}\right)^2$  *wow!!*  
*constant, doesn't depend on T*  
*dimensionality of X !!*

©Carlos Guestrin 2005-2013

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

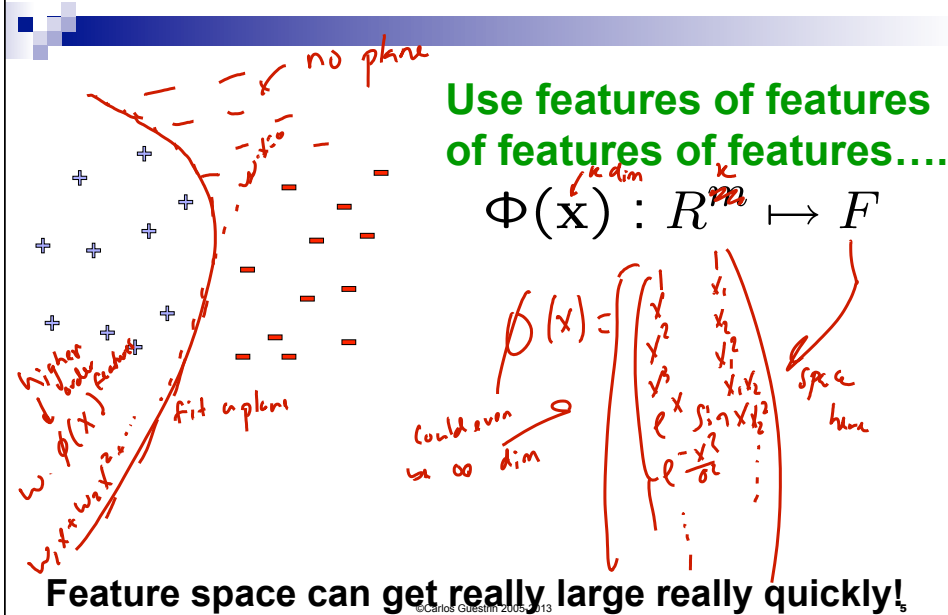


*we need features that make data as linearly separable as possible*

©Carlos Guestrin 2005-2013

4

What if the data is not linearly separable?

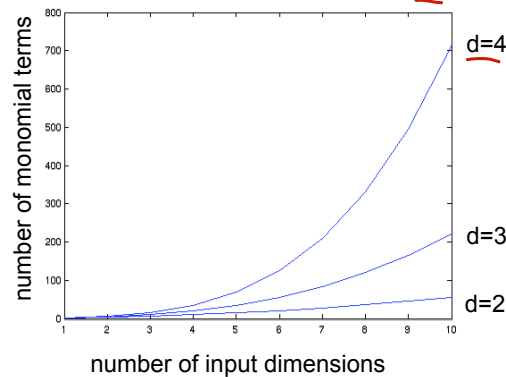


## Higher order polynomials

$$\text{num. terms} = \binom{d+m-1}{d} = \frac{(d+m-1)!}{d!(m-1)!}$$

dim of  $\phi(x) = \text{dim of } w$

$m$  – input features  
 $d$  – degree of polynomial



Even though  
dims of  $\phi(x)$   
are huge, AI model  
very ~~not~~ quickly  
grows fast!  
 $d = 6, m = 100$   
about 1.6 billion terms

# Perceptron Revisited

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

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

- Mistake at time  $t$ :  $w^{(t+1)} \leftarrow w^{(t)} + y^{(t)} x^{(t)}$

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

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

$$w^{(t)} = \sum_{j \in M^{(t)}} y^{(j)} x^{(j)}$$

- Prediction rule now:

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

- When using high dimensional features:

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

when can you compute this efficiently

for simplicity  $\phi(0) = 0$   
 how similar  $x$  is to my old mistakes  
 classification rule only depends on dot prod. between points

©Carlos Guestrin 2005-2013

7

# Dot-product of polynomials

$\Phi(u) \cdot \Phi(v)$  = polynomials of degree exactly  $d$

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

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

proof by base case (one step of induction)

for poly of degree exactly  $d$

$$\phi(u) \cdot \phi(v) = (u \cdot v)^d$$

compute dot product extremely efficiently

kernel trick in general

©Carlos Guestrin 2005-2013

8

## Finally the Kernel Trick!!! (Kernelized Perceptron)

- Every time you make a mistake, remember  $(\mathbf{x}^{(t)}, y^{(t)})$   
 $\hookrightarrow$  keep indices  $M(t)$  mistakes up to time  $t$

- Kernelized Perceptron prediction for  $\mathbf{x}$ :

prediction

$$\text{sign}(\mathbf{w}^{(t)} \cdot \phi(\mathbf{x})) = \sum_{j \in M^{(t)}} y^{(j)} \phi(\mathbf{x}^{(j)}) \cdot \phi(\mathbf{x})$$

perceptron alg with kernels  
 $\rightarrow$  keep mistakes in memory  
 $\rightarrow$  compute prediction if mistake, save it

$$= \sum_{j \in M^{(t)}} y^{(j)} k(\mathbf{x}^{(j)}, \mathbf{x})$$

kernel function  
 $k(u, v) = \phi(u) \cdot \phi(v)$

©Carlos Guestrin 2005-2013

9

## Polynomial kernels

- All monomials of degree  $d$  in  $O(d)$  operations:  
 $\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d = \text{polynomials of degree exactly } d$

- How about all monomials of degree up to  $d$ ?

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

□ Better solution:  $d=2, 1$   
 $(\mathbf{u} \cdot \mathbf{v})^2 + (\mathbf{u} \cdot \mathbf{v})^1 + (\mathbf{u} \cdot \mathbf{v})^0 = (\mathbf{u} \cdot \mathbf{v} + 1)^2$

for polynomials of degree  $d$   
 $\phi(\mathbf{u}) \cdot \phi(\mathbf{v}) = k(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$

©Carlos Guestrin 2005-2013

10

## Common kernels

MANY OTHERS

- 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

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

o strings  
o graph  
o ...

$u_1, u_2, \dots$  attributes are fixed  
 $\phi(\cdot)$  projects them into  
high dim space  
equivalent to  $\phi(\mathbf{u}) \cdot \phi(\mathbf{v})$   
where  $\phi(\mathbf{u})$  is infinite  
dimensional  
dot prod. in  
infinite space

©Carlos Guestrin 2005-2013

11

## 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

©Carlos Guestrin 2005-2013

12

$\approx$  perception + regularization

# Support Vector Machines

Machine Learning – CSE446  
Carlos Guestrin  
University of Washington  
May 3, 2013  
©Carlos Guestrin 2005-2013

13

## Linear classifiers – Which line is better?

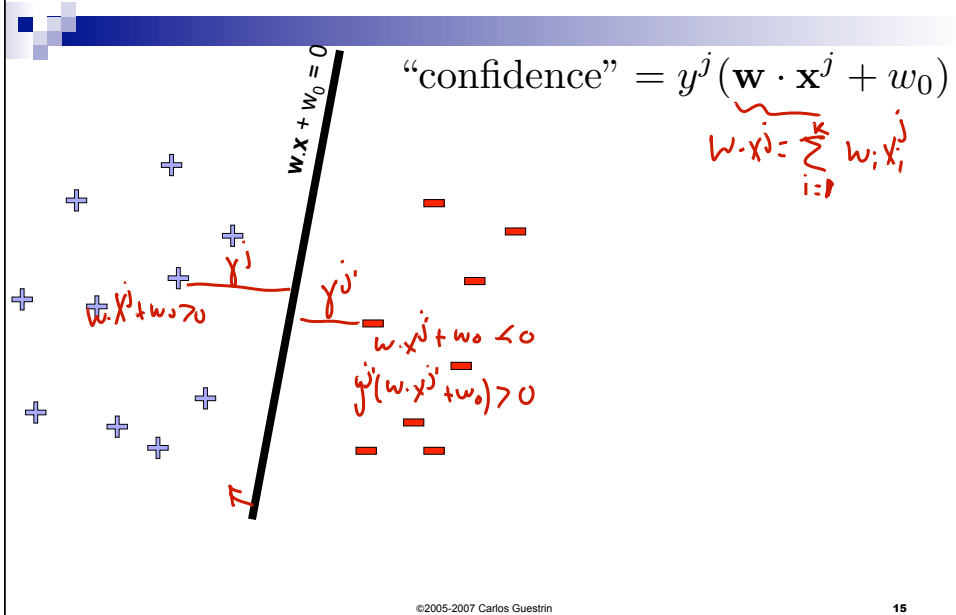
furthest away from nearby points  $\equiv$  margin maximizing

$w \cdot x + w_0 = 0$

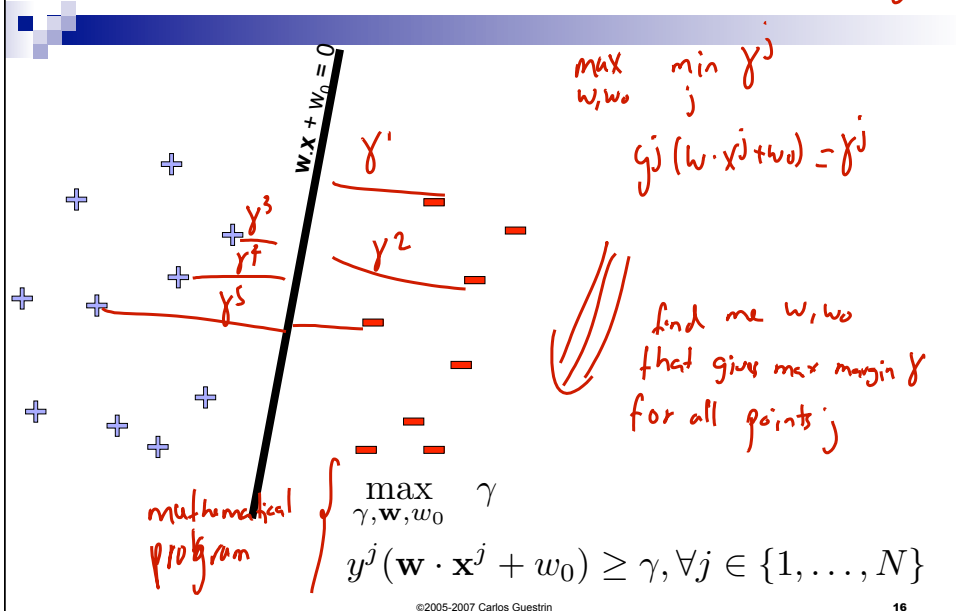
©2005-2007 Carlos Guestrin

14

# Pick the one with the largest margin!

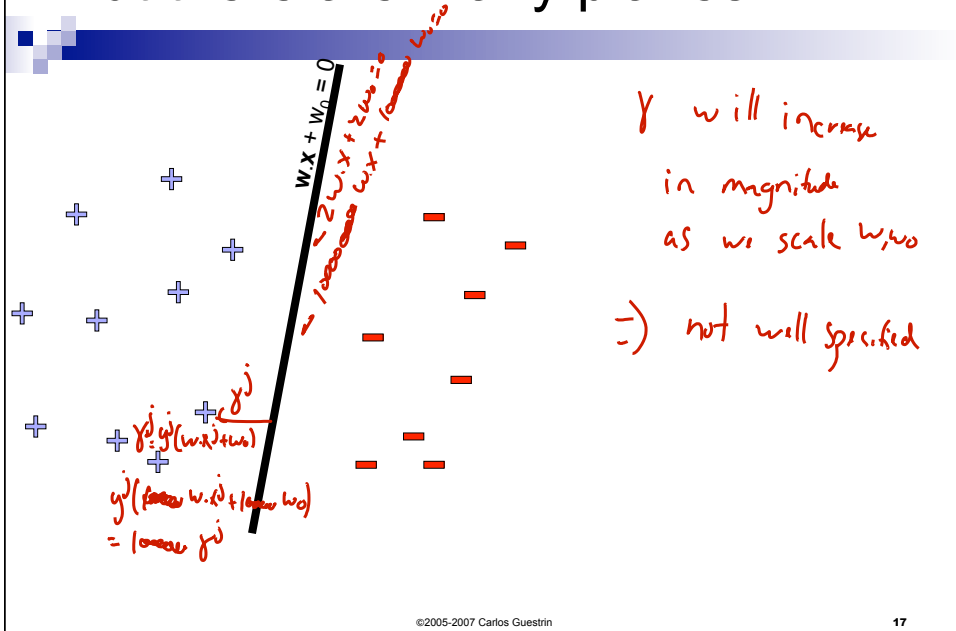


# Maximize the margin

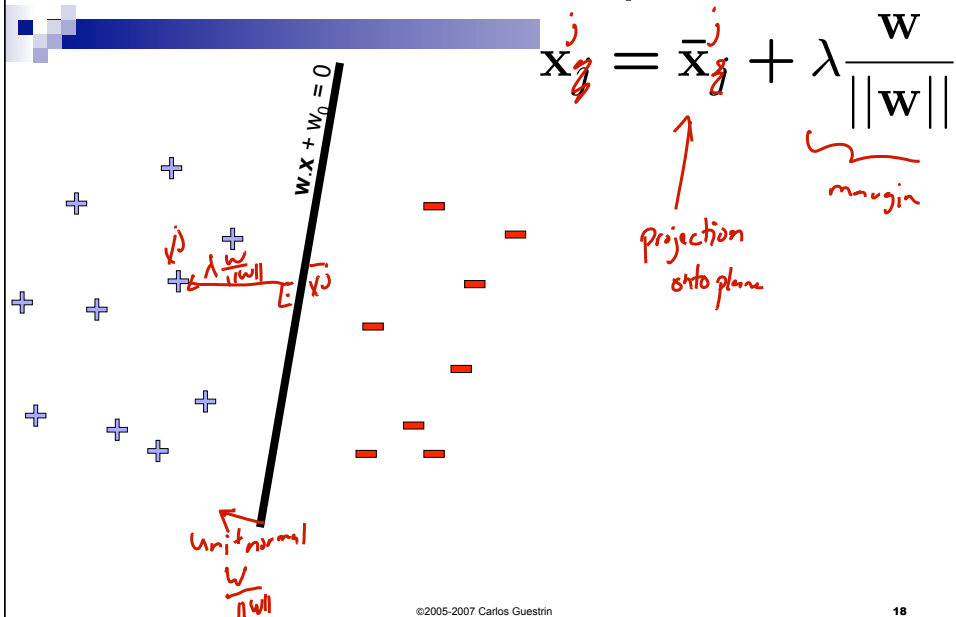


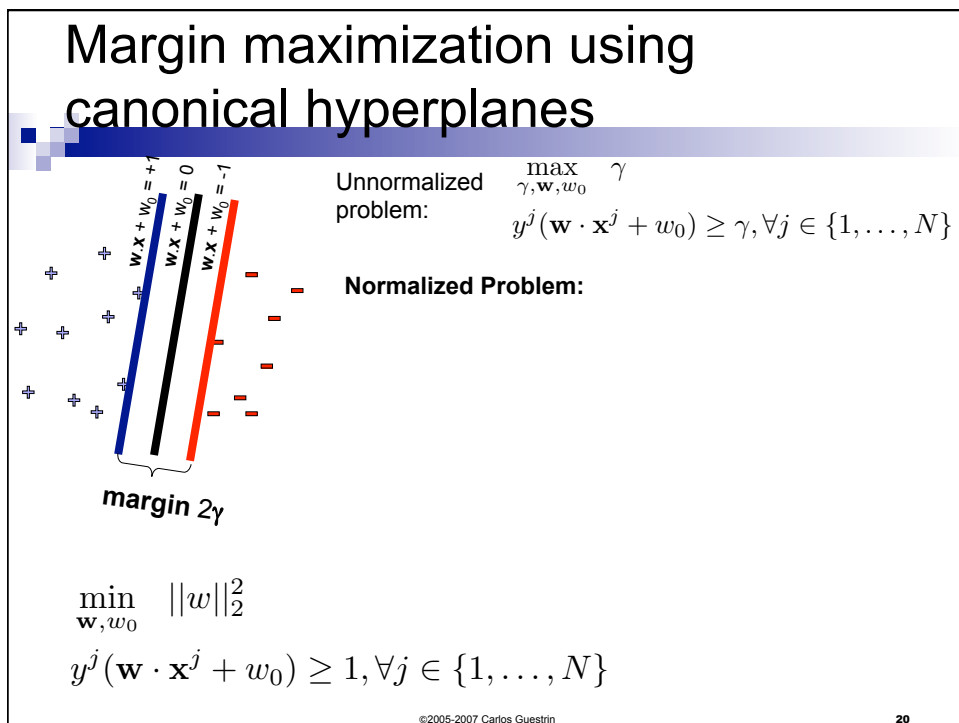
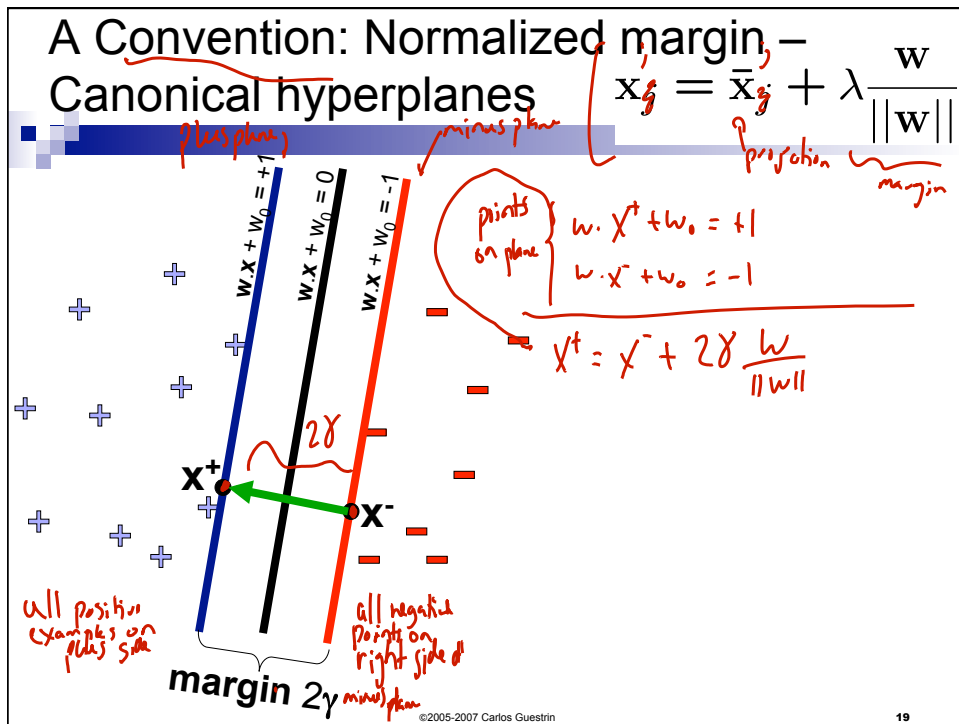


## But there are many planes...

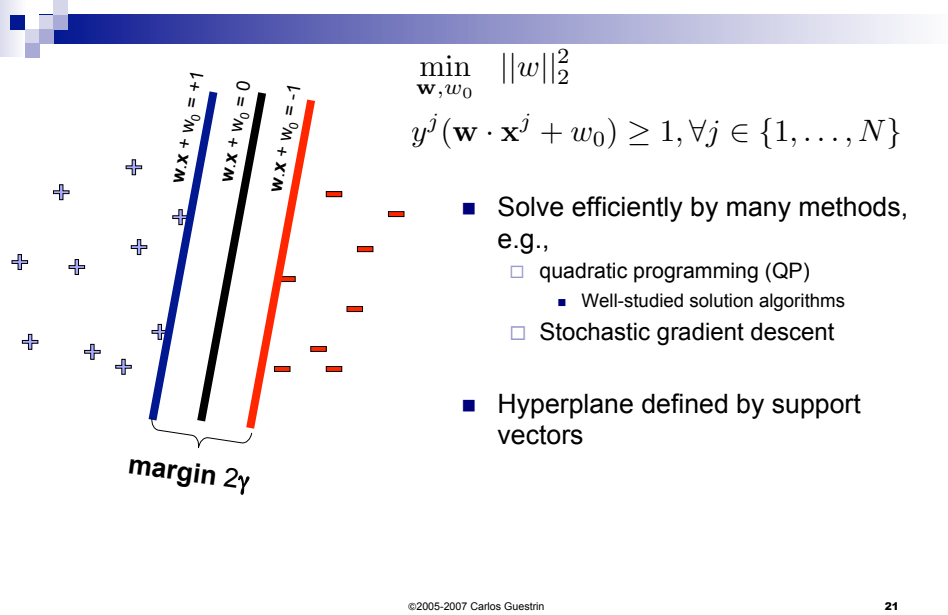


## Review: Normal to a plane





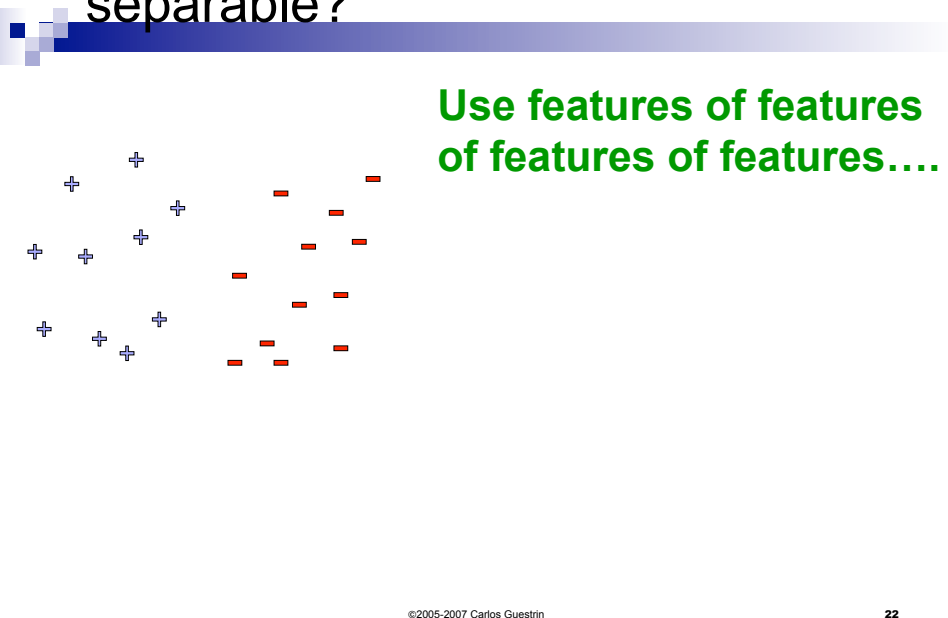
# Support vector machines (SVMs)



©2005-2007 Carlos Guestrin

21

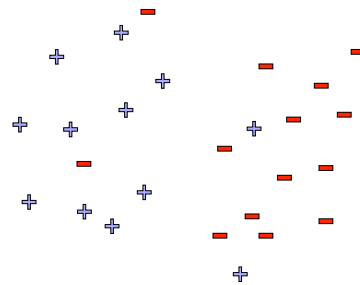
## What if the data is not linearly separable?



©2005-2007 Carlos Guestrin

22

## What if the data is still not linearly separable?



$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

$$y^j(\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1, \forall j$$

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

©2005-2007 Carlos Guestrin

23

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



- Perceptron was minimizing the hinge loss:

$$\sum_{j=1}^N (-y^j(\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

- SVMs minimizes the regularized hinge loss!!

$$\|\mathbf{w}\|_2^2 + C \sum_{j=1}^N (1 - y^j(\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

©Carlos Guestrin 2005-2013

24

## Stochastic Gradient Descent for SVMs

- Perceptron minimization:

$$\sum_{j=1}^N (-y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

- SGD for Perceptron:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} [y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0] y^{(t)} \mathbf{x}^{(t)}$$

- SVMs minimization:

$$\|\mathbf{w}\|_2^2 + C \sum_{j=1}^N (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

- SGD for SVMs:

©Carlos Guestrin 2005-2013

25

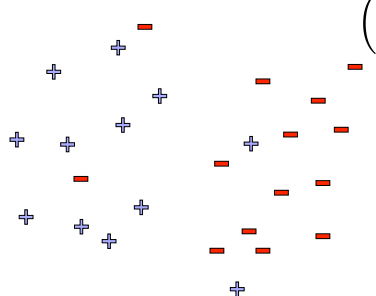
## What you need to know

- 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

©2005-2007 Carlos Guestrin

26

## Slack variables – Hinge loss



$$\text{minimize}_{\mathbf{w}, b} \quad \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j$$

$$\xi_j \geq 0, \quad \forall j$$

- If margin  $\geq 1$ , don't care
- If margin  $< 1$ , pay linear penalty

©2005-2007 Carlos Guestrin 27

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

**SVM:**

$$\text{minimize}_{\mathbf{w}, b} \quad \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j$$

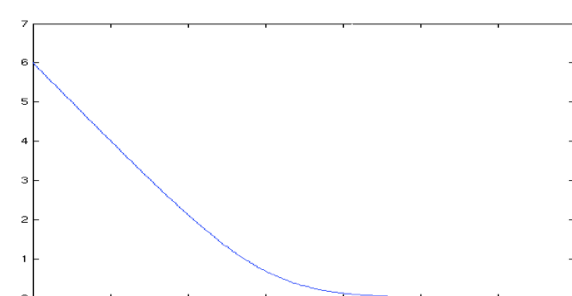
$$\xi_j \geq 0, \quad \forall j$$

**Logistic regression:**

$$P(Y = 1 | x, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}}$$

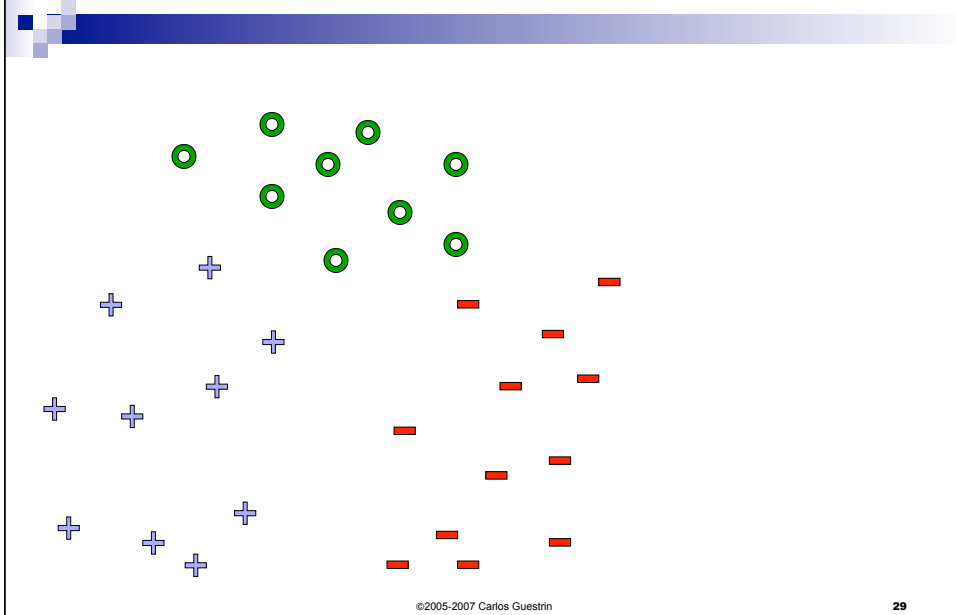
**Log loss:**

$$-\ln P(Y = 1 | x, \mathbf{w}) = \ln(1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)})$$



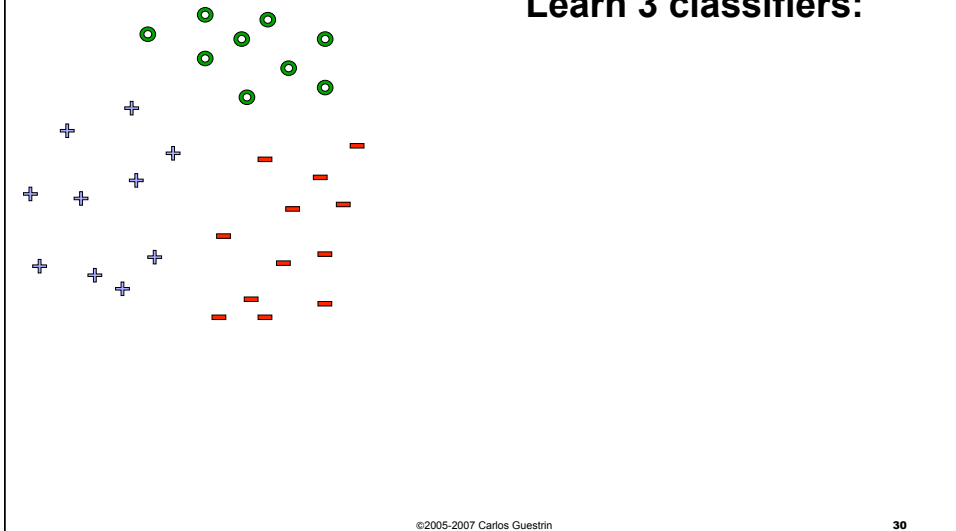
©2005-2007 Carlos Guestrin 28

## What about multiple classes?



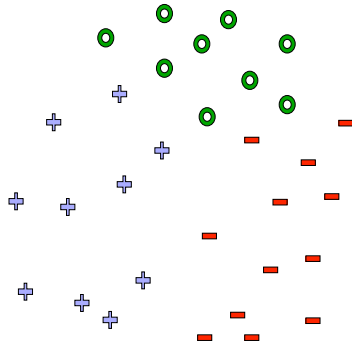
## One against All

**Learn 3 classifiers:**



## Learn 1 classifier: Multiclass SVM

Simultaneously learn 3 sets of weights



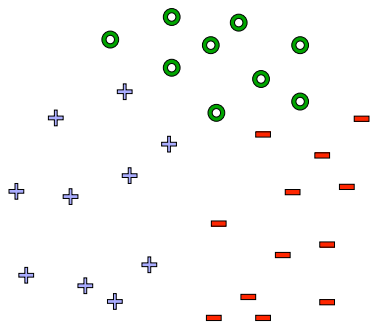
$$\mathbf{w}^{(y_j)} \cdot \mathbf{x}_j + b^{(y_j)} \geq \mathbf{w}^{(y')} \cdot \mathbf{x}_j + b^{(y')} + 1, \quad \forall y' \neq y_j, \quad \forall j$$

©2005-2007 Carlos Guestrin

31

## Learn 1 classifier: Multiclass SVM

$$\begin{aligned} &\text{minimize}_{\mathbf{w}, b} \quad \sum_y \mathbf{w}^{(y)} \cdot \mathbf{w}^{(y)} + C \sum_j \xi_j \\ &\mathbf{w}^{(y_j)} \cdot \mathbf{x}_j + b^{(y_j)} \geq \mathbf{w}^{(y')} \cdot \mathbf{x}_j + b^{(y')} + 1 - \xi_j, \quad \forall y' \neq y_j, \quad \forall j \\ &\xi_j \geq 0, \quad \forall j \end{aligned}$$



©2005-2007 Carlos Guestrin

32