

~perception + regularization

Support Vector Machines

Machine Learning – CSE446

Carlos Guestrin

University of Washington

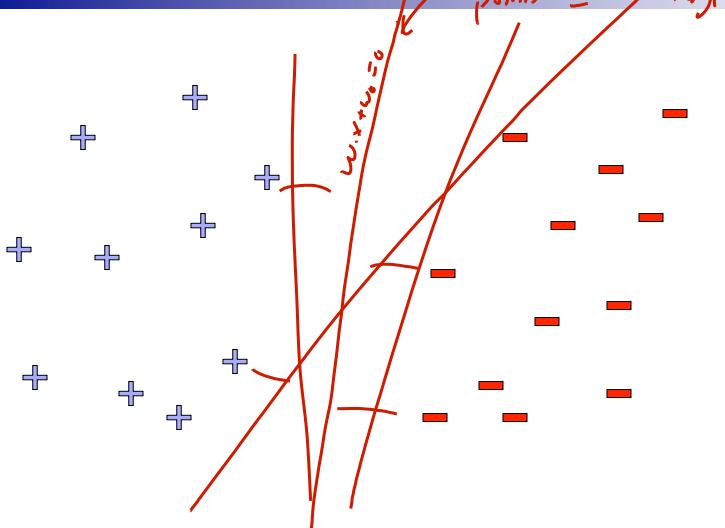
May 6, 2013

©Carlos Guestrin 2005-2013

1

Linear classifiers – Which line is better?

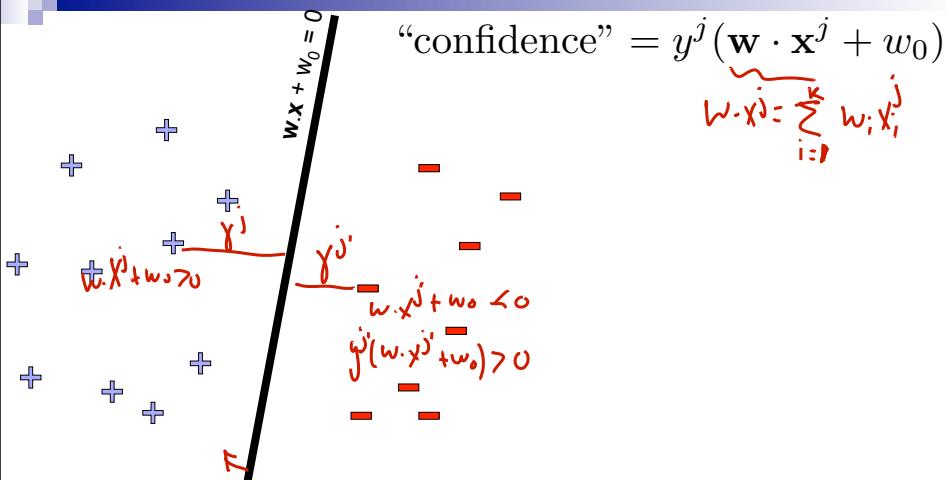
furthest away from nearby points \equiv margin maximizing



©Carlos Guestrin 2005-2013

2

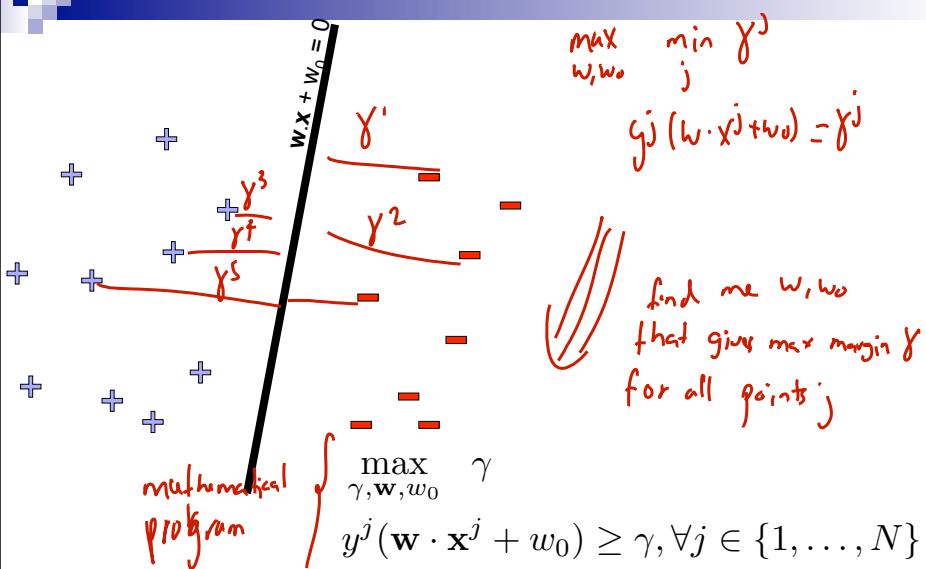
Pick the one with the largest margin!



©Carlos Guestrin 2005-2013

3

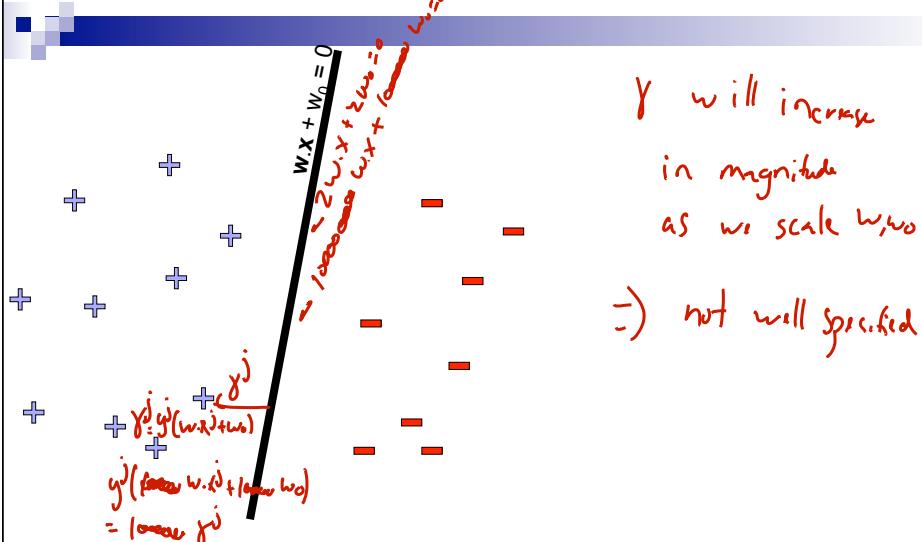
Maximize the margin *maximize worst case margin*



©Carlos Guestrin 2005-2013

4

But there are many planes...



©Carlos Guestrin 2005-2013

5

Review: Normal to a plane

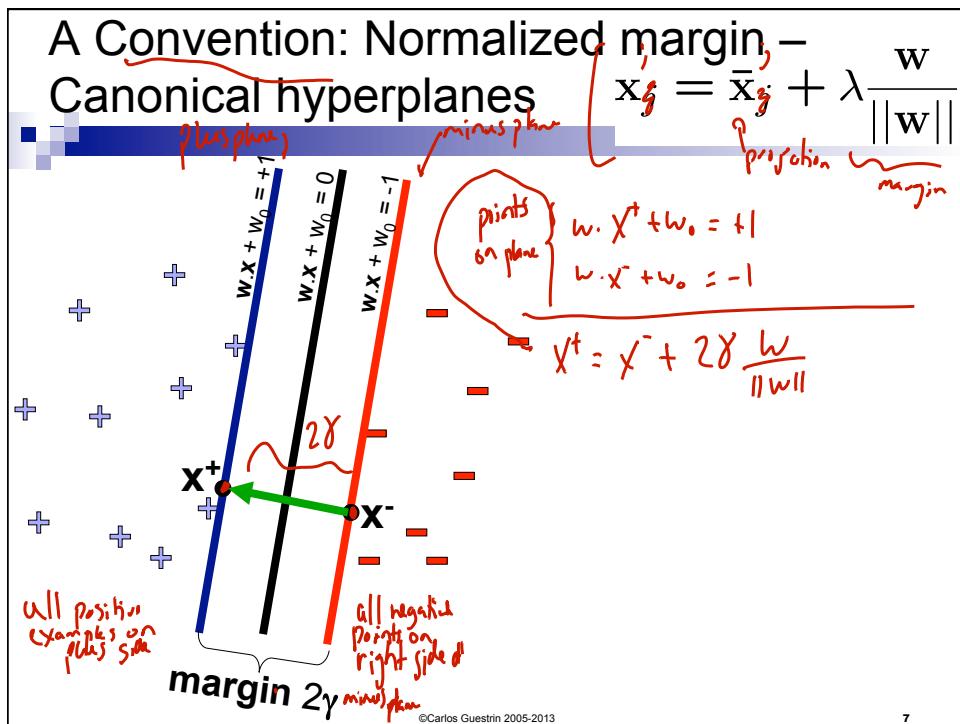
$$\mathbf{x}_j^j = \bar{\mathbf{x}}_j^j + \lambda \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

projection
onto plane

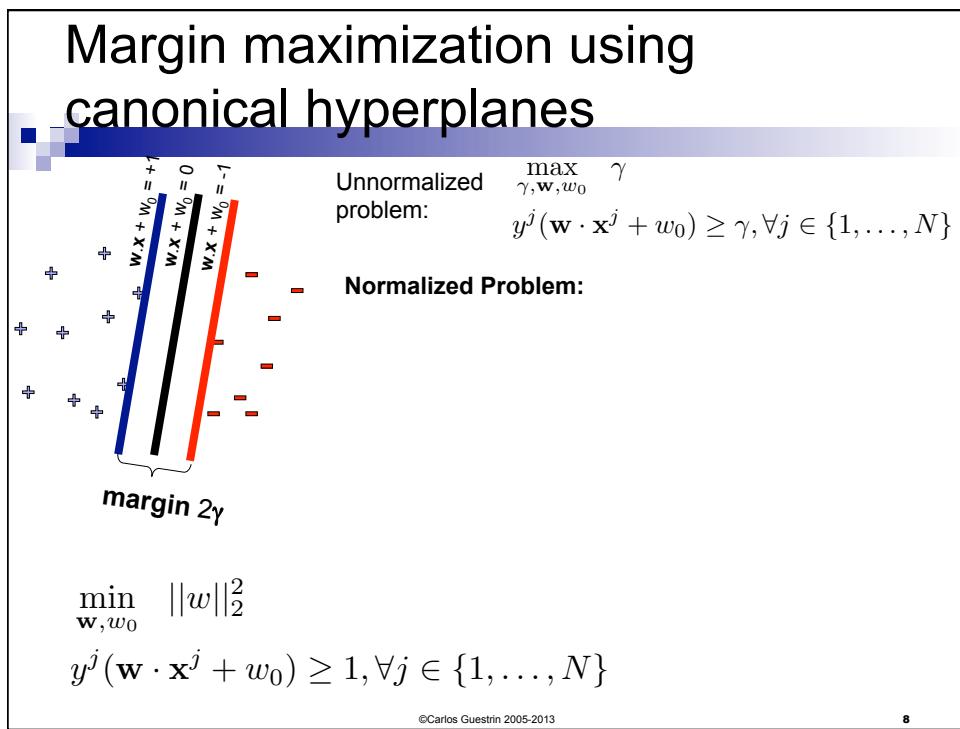
mugia

©Carlos Guestrin 2005-2013

6

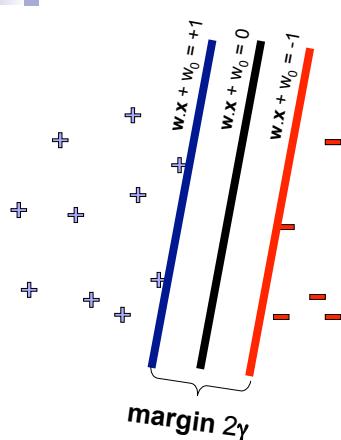


7



8

Support vector machines (SVMs)



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

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

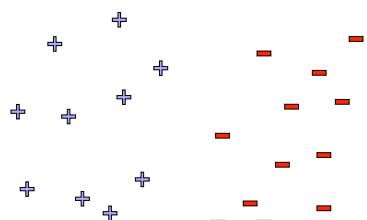
- Solve efficiently by many methods, e.g.,
 - quadratic programming (QP)
 - Well-studied solution algorithms
 - Stochastic gradient descent
- Hyperplane defined by support vectors

©Carlos Guestrin 2005-2013

9

What if the data is not linearly separable?

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



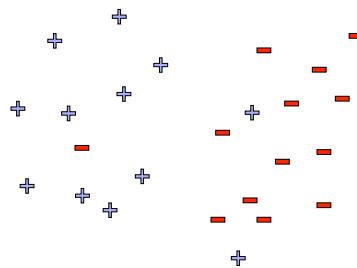
©Carlos Guestrin 2005-2013

10

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 $\|\mathbf{w}\|$:

©Carlos Guestrin 2005-2013

11

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

12

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

13

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

©Carlos Guestrin 2005-2013

14

Big Picture

Machine Learning – CSE446

Carlos Guestrin

University of Washington

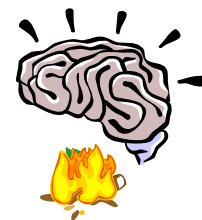
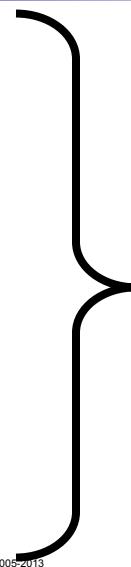
May 6, 2013

©Carlos Guestrin 2005-2013

15

What you have learned thus far

- Learning is function approximation
- Point estimation
- Regression
- LASSO
- Logistic regression
- Bias-Variance tradeoff
- Regularization
- Decision trees
- Cross validation
- Boosting
- Instance-based learning
- Online learning
- Perceptron
- SVMs
- Kernel trick



©Carlos Guestrin 2005-2013

16

Review material in terms of...

- Types of learning problems
- Hypothesis spaces
- Loss functions
- Optimization algorithms

©Carlos Guestrin 2005-2013

17

ML Pipeline

Attributes/ Observations	Features/ Basis Functions	Task	Hypothesis Class/ Model	Algorithm/ Optimization Method
-----------------------------	------------------------------	------	----------------------------	-----------------------------------

©Carlos Guestrin 2005-2013

18

Learning Task/Measuring Error

TASK	LOSS FUNCTIONS
Regression	
Classification	
Density Estimation	

©Carlos Guestrin 2005-2013

19

Hypothesis Classes & Decision Boundaries

Simple Linear Model

Linear Model with
Higher-Order Features or Kernels

Nearest Neighbors

Boosting

©Carlos Guestrin 2005-2013

20

The Power of Regularization

Overfitting
Bias/Variance Tradeoff

Regularization

©Carlos Guestrin 2005-2013

21

Your Midterm...

- Content: Everything up to today...
- Only 50mins, so arrive early and settle down quickly
- “Open book”
 - Textbook, Course notes, Personal notes
- No:
 - Computer, phone, other materials,...
- The exam:
 - Covers key concepts and ideas, work on understanding the big picture, and differences between methods

©Carlos Guestrin 2005-2013

22