

perception + regularization

Support Vector Machines

Machine Learning – CSE446

Carlos Guestrin

University of Washington

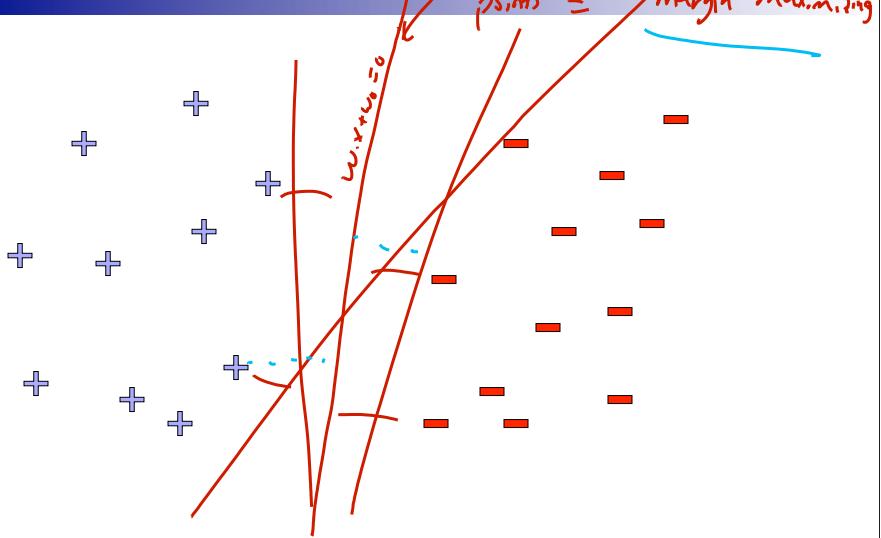
May 6, 2013

©Carlos Guestrin 2005-2013

1

Linear classifiers – Which line is better?

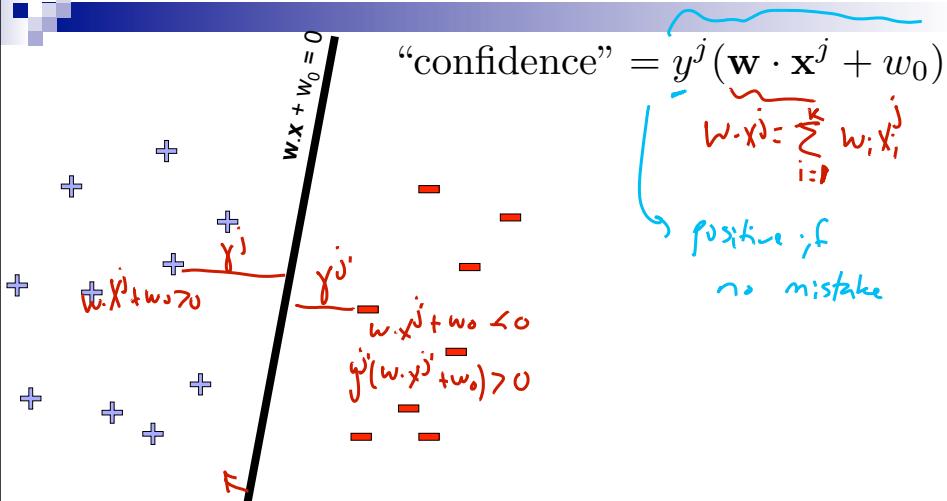
furthest away from nearby points \equiv margin maximizing



©Carlos Guestrin 2005-2013

2

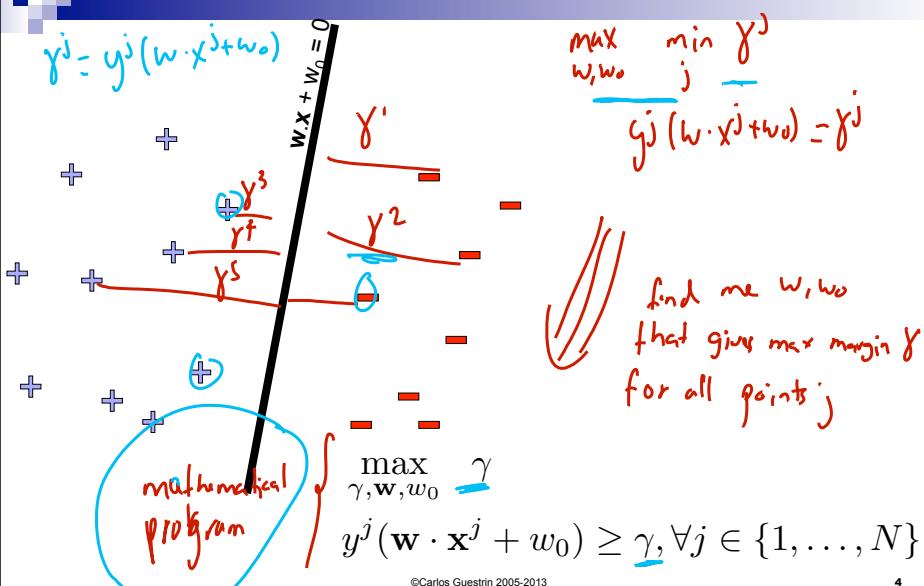
Pick the one with the largest margin!



©Carlos Guestrin 2005-2013

3

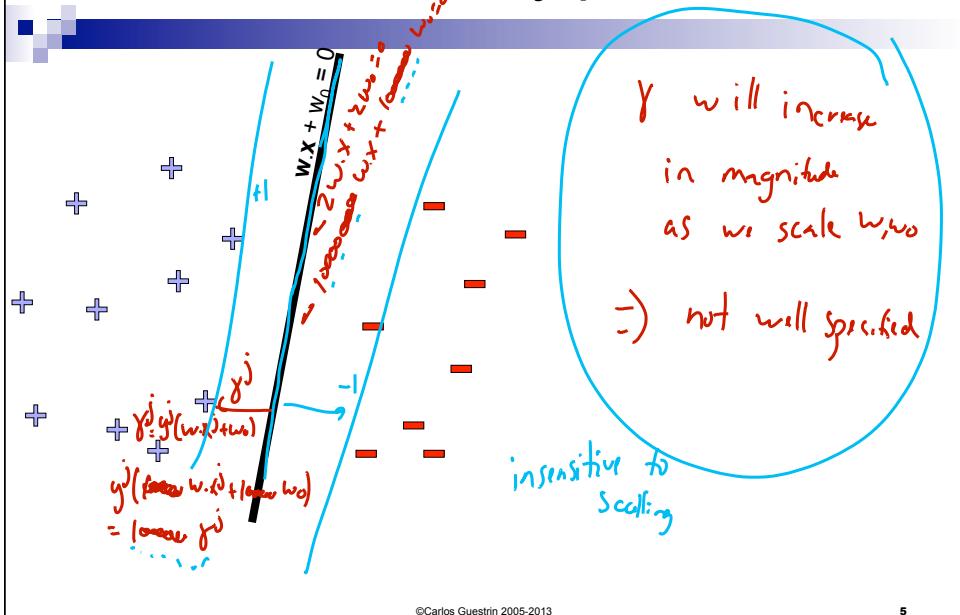
Maximize the margin *maximize worst case margin*



©Carlos Guestrin 2005-2013

4

But there are many planes...



©Carlos Guestrin 2005-2013

5

Review: Normal to a plane

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

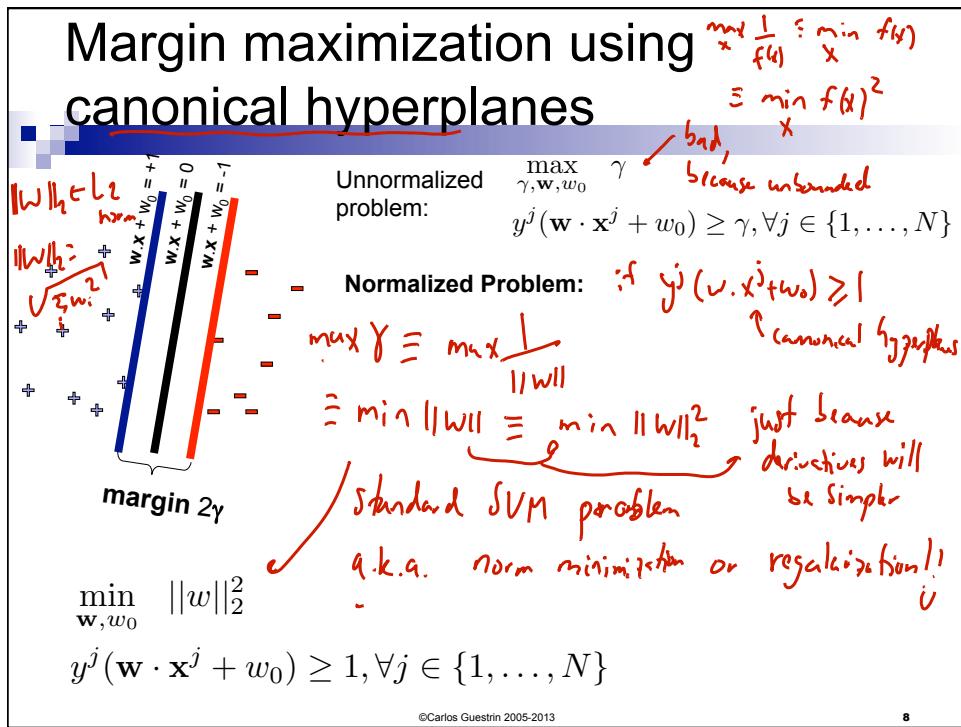
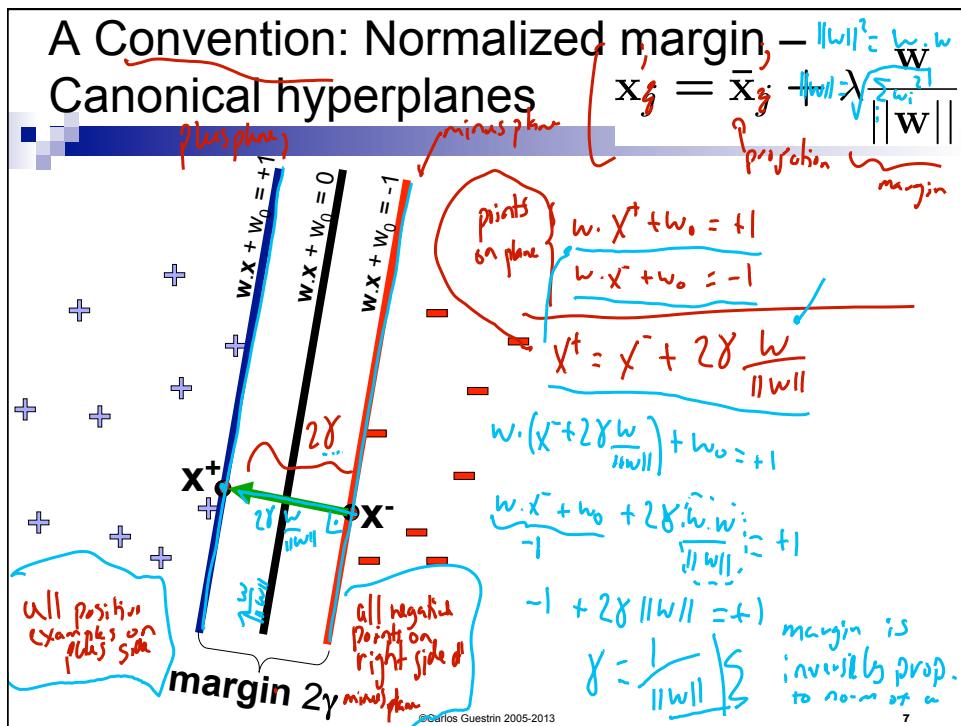
projection onto plane

margin

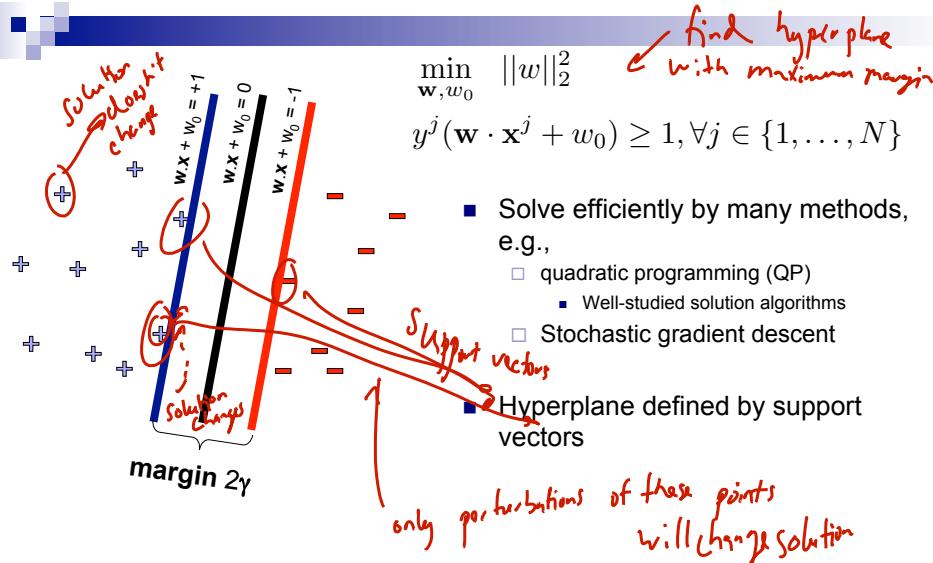
Unit normal
 $\frac{\mathbf{w}}{\|\mathbf{w}\|}$

©Carlos Guestrin 2005-2013

6



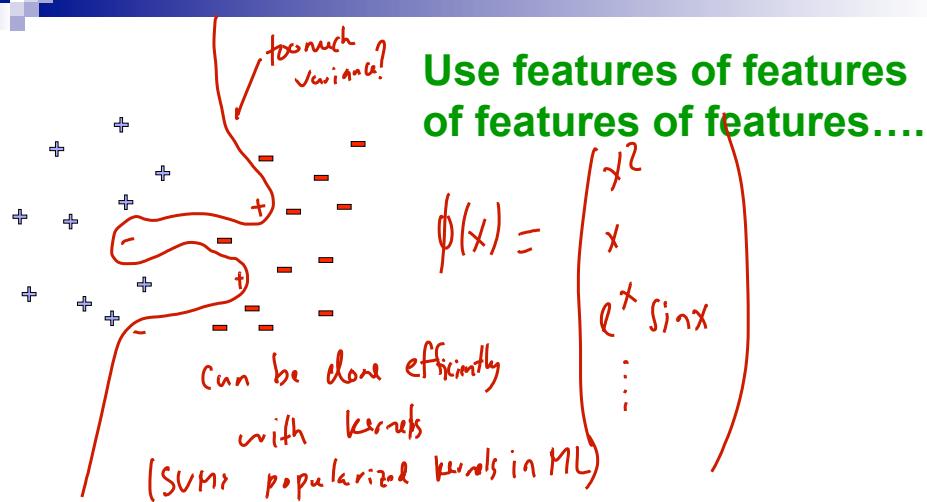
Support vector machines (SVMs)



©Carlos Guestrin 2005-2013

9

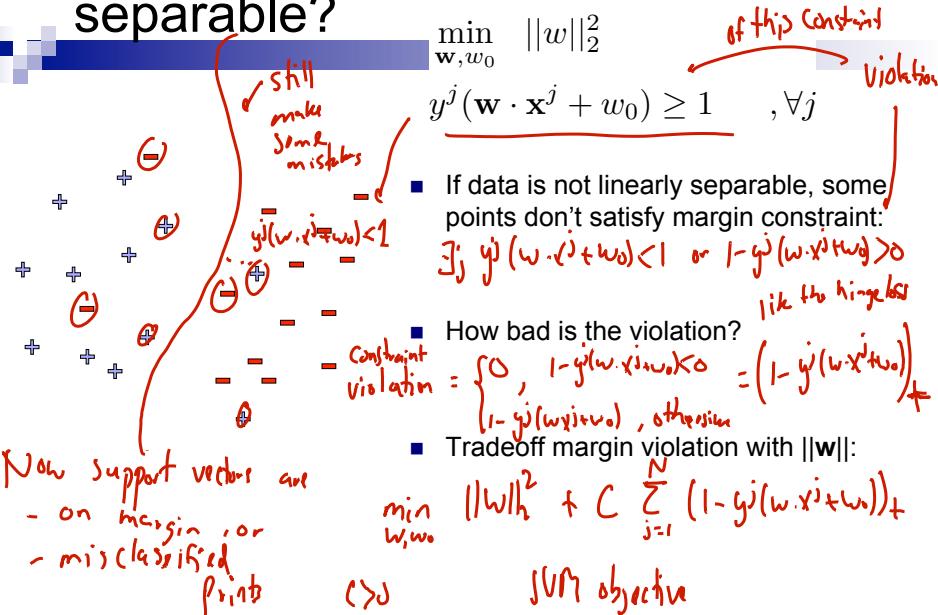
What if the data is not linearly separable?



©Carlos Guestrin 2005-2013

10

What if the data is still not linearly separable?



©Carlos Guestrin 2005-2013

11

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

- Perceptron was minimizing the hinge loss:

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

both hinge loss

- SVMs minimizes the regularized hinge loss!!

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

just convention

extra regularization!!

©Carlos Guestrin 2005-2013

12

$$\nabla_w \|\mathbf{w}\|_2^2 = \nabla_w \mathbf{w} \cdot \mathbf{w} = 2\mathbf{w} \leftarrow \text{direction of increase}$$

but we go forward - $\nabla_w \|\mathbf{w}\|_2^2$

Stochastic Gradient Descent for SVMs

- Perceptron minimization:

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

- SGD for Perceptron:

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

make mistake
update weights
step size = 1

- 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:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \left(C \mathbb{1}_{[(1 - g^{(t)}(\mathbf{w} \cdot \mathbf{x}^{(t)} + w_0)) > 0]} g^{(t)} \mathbf{x}^{(t)} - 2\mathbf{w}^{(t)} \right)$$

skip size needs to agree with other ways
pick C, η by cross validation

©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

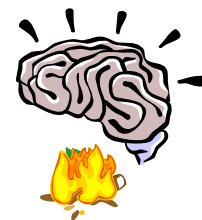
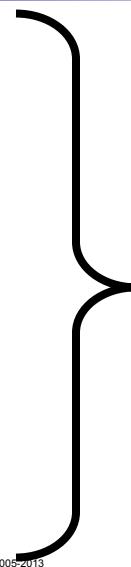
May 6, 2013

©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
X	$h_i(x)$	Regression	Linear models	Optimize a loss function
age	$\phi(x)$	$X \text{ or } \phi(x) \rightarrow \mathbb{R}$	$w \cdot x$	- Gradient
gender	$\begin{cases} \text{age}^2 \\ \text{age} \end{cases}$	Classification	NN, DTs	- Set gradient = 0 closed form
wage	$\begin{cases} \text{age}^2 \\ \text{age} \\ \vdots \\ \vdots \end{cases}$	$X \text{ or } \phi(x) \rightarrow \{1, \dots, r\}$	Boosting	- Stochastic Coordinate descent
:		Density estimation		- SGD
		$X \text{ or } \phi(x) \rightarrow \text{probability}$		- Search decision trees
		Logistic regression		

©Carlos Guestrin 2005-2013

18

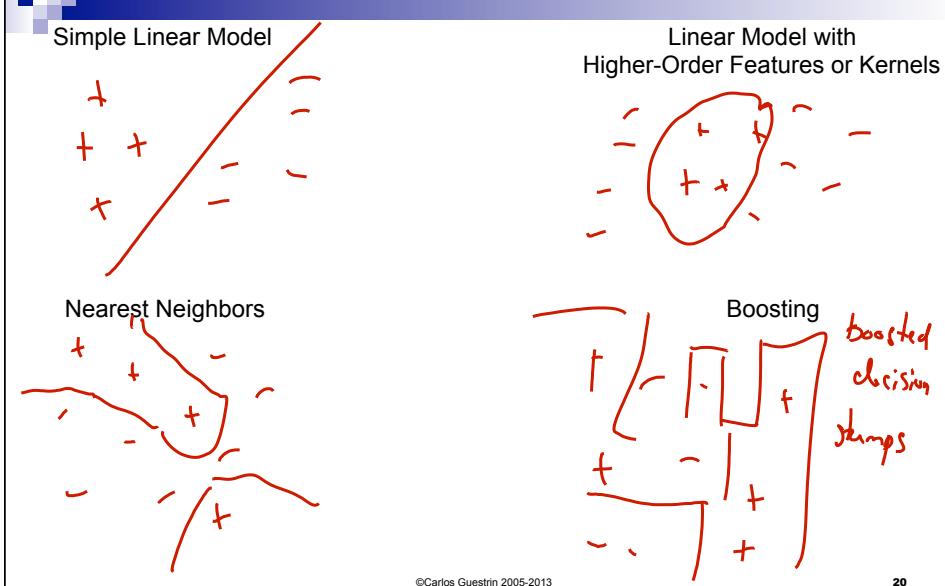
Learning Task/Measuring Error

TASK	LOSS FUNCTIONS
Regression	Squared error
Classification	log loss for logistic regression hinge loss
Density Estimation	log loss for LR

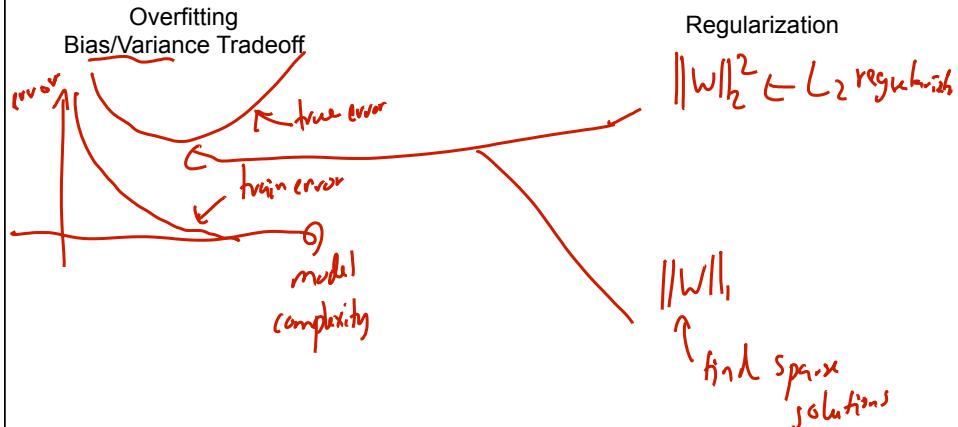
©Carlos Guestrin 2005-2013

19

Hypothesis Classes & Decision Boundaries



The Power of 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