

# Instance-based Learning

Machine Learning – 10701/15781

Carlos Guestrin

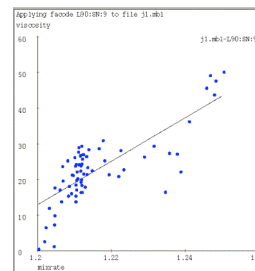
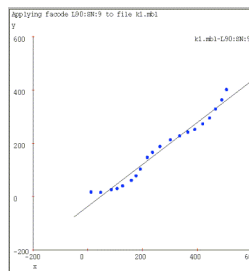
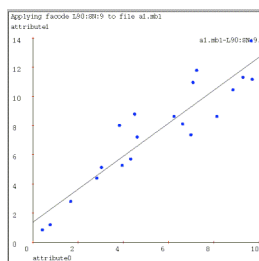
Carnegie Mellon University

February 19<sup>th</sup>, 2007

©2005-2007 Carlos Guestrin

1

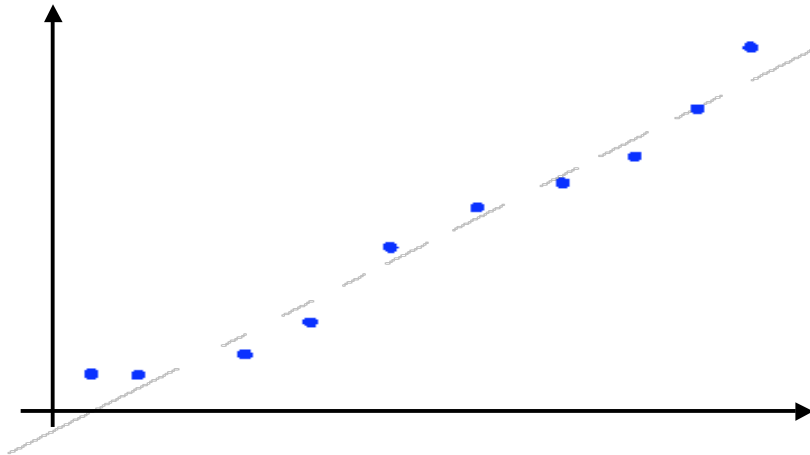
## Why not just use Linear Regression?



©2005-2007 Carlos Guestrin

2

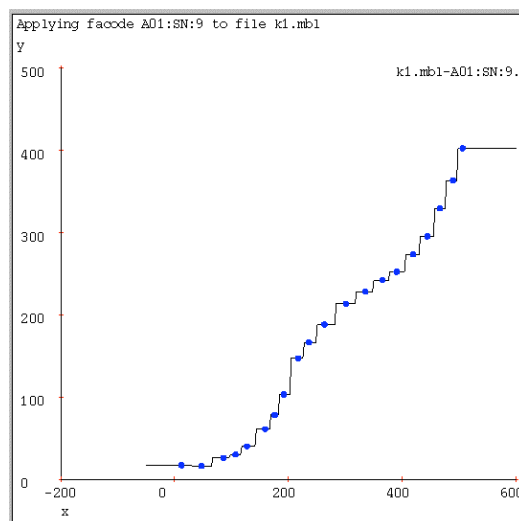
## Using data to predict new data



©2005-2007 Carlos Guestrin

3

## Nearest neighbor



©2005-2007 Carlos Guestrin

4

# Univariate 1-Nearest Neighbor

Given datapoints  $(x_1, y_1) (x_2, y_2) \dots (x_N, y_N)$ , where we assume  $y_i = f(x_i)$  for some unknown function  $f$ .

Given query point  $x_q$ , your job is to predict  $\hat{y} \approx f(x_q)$

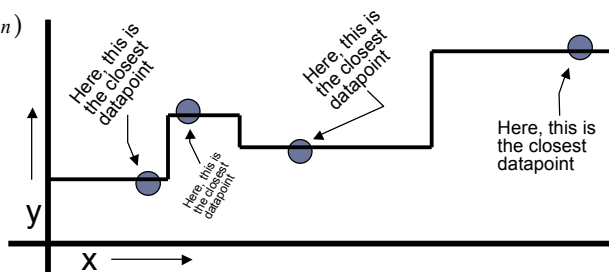
Nearest Neighbor:

1. Find the closest  $x_i$  in our set of datapoints

$$i(nn) = \underset{i}{\operatorname{argmin}} |x_i - x_q|$$

2. Predict  $\hat{y} = y_{i(nn)}$

Here's a dataset with one input, one output and four datapoints.



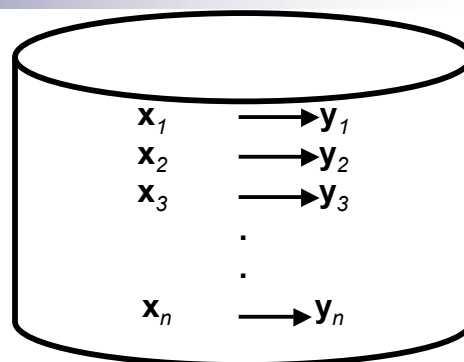
©2005-2007 Carlos Guestrin

5

## 1-Nearest Neighbor is an example of.... Instance-based learning

A function approximator that has been around since about 1910.

To make a prediction, search database for similar datapoints, and fit with the local points.



**Four things make a memory based learner:**

- A distance metric
- How many nearby neighbors to look at?
- A weighting function (optional)
- How to fit with the local points?

©2005-2007 Carlos Guestrin

6

# 1-Nearest Neighbor

Four things make a memory based learner:

1. *A distance metric*  
**Euclidian (and many more)**
2. *How many nearby neighbors to look at?*  
**One**
3. *A weighting function (optional)*  
**Unused**
4. *How to fit with the local points?*  
**Just predict the same output as the nearest neighbor.**

©2005-2007 Carlos Guestrin

7

## Multivariate 1-NN examples

Regression

Classification

©2005-2007 Carlos Guestrin

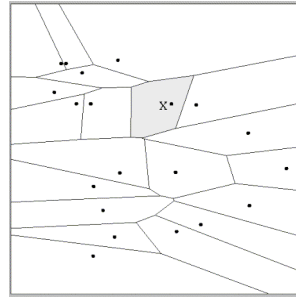
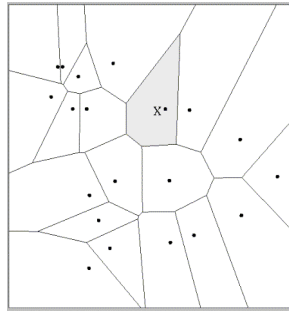
8

# Multivariate distance metrics

Suppose the input vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  are two dimensional:

$\mathbf{x}_1 = (x_{11}, x_{12}), \mathbf{x}_2 = (x_{21}, x_{22}), \dots, \mathbf{x}_N = (x_{N1}, x_{N2})$ .

One can draw the nearest-neighbor regions in input space.



$$Dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 \quad Dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (3x_{i2} - 3x_{j2})^2$$

The relative scalings in the distance metric affect region shapes

©2005-2007 Carlos Guestrin

9

# Euclidean distance metric

Or equivalently,

$$D(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_i \sigma_i^2 (x_i - x'_i)^2}$$

where

$$D(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \Sigma (\mathbf{x} - \mathbf{x}')}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & \sigma_N^2 \end{bmatrix}$$

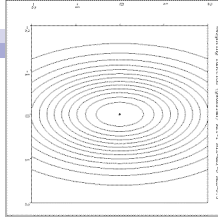
Other Metrics...

- Mahalanobis, Rank-based, Correlation-based,...

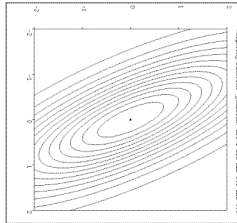
©2005-2007 Carlos Guestrin

10

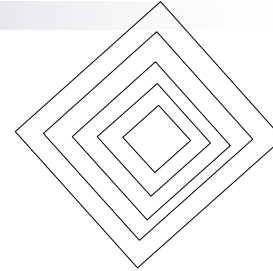
## Notable distance metrics (and their level sets)



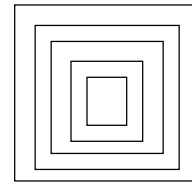
**Scaled Euclidian ( $L_2$ )**



**Mahalanobis**  
(here,  $\Sigma$  on the previous  
slide is not necessarily  
diagonal, but is symmetric)



**$L_1$  norm (absolute)**



**$L_\infty$  (max) norm**

11

©2005-2007 Carlos Guestrin

## Consistency of 1-NN

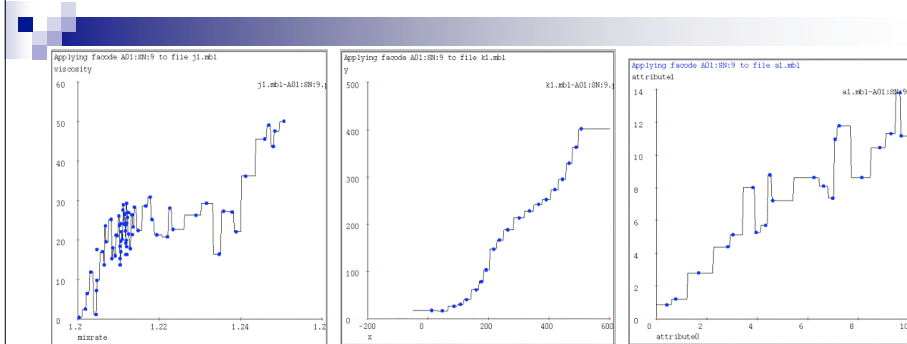
- Consider an estimator  $f_n$  trained on  $n$  examples
    - e.g., 1-NN, neural nets, regression,...
  - Estimator is *consistent* if true error goes to zero as amount of data increases
    - e.g., for no noise data, consistent if:
- $$\lim_{n \rightarrow \infty} MSE(f_n) = 0$$
- Regression is not consistent!
    - Representation bias
  - **1-NN is consistent** (under some mild fineprint)

**What about variance???**

©2005-2007 Carlos Guestrin

12

# 1-NN overfits?



©2005-2007 Carlos Guestrin

13

## k-Nearest Neighbor

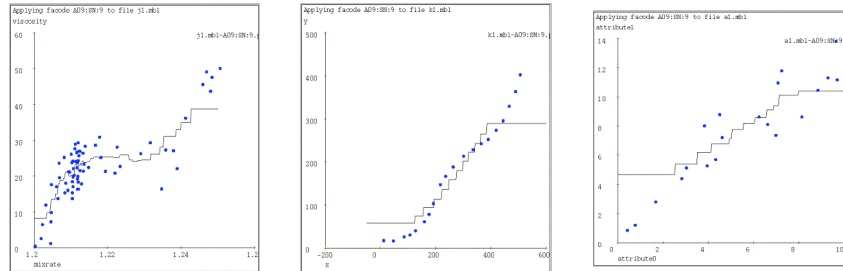
Four things make a memory based learner:

1. *A distance metric*  
**Euclidian (and many more)**
2. *How many nearby neighbors to look at?*  
**k**
1. *A weighting function (optional)*  
**Unused**
2. *How to fit with the local points?*  
**Just predict the average output among the k nearest neighbors.**

©2005-2007 Carlos Guestrin

14

## k-Nearest Neighbor (here k=9)



**K-nearest neighbor for function fitting smooths away noise, but there are clear deficiencies.**

What can we do about all the discontinuities that k-NN gives us?

©2005-2007 Carlos Guestrin

15

## Weighted k-NNs

- Neighbors are not all the same

©2005-2007 Carlos Guestrin

16



# Kernel regression

Four things make a memory based learner:

1. A distance metric  
**Euclidian (and many more)**
2. How many nearby neighbors to look at?  
**All of them**
3. A weighting function (optional)  
 $w_i = \exp(-D(x_i, \text{query})^2 / K_w^2)$

Nearby points to the query are weighted strongly, far points weakly. The  $K_w$  parameter is the **Kernel Width**. Very important.

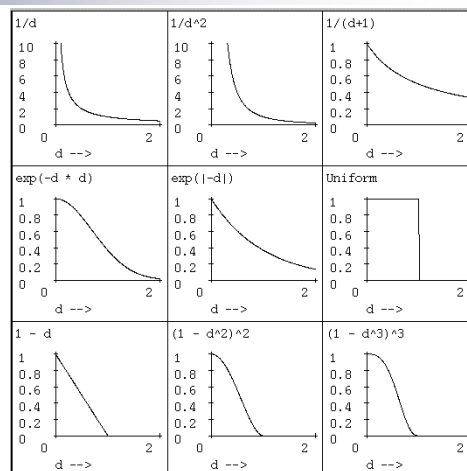
4. How to fit with the local points?  
**Predict the weighted average of the outputs:**  
 $\text{predict} = \sum w_i y_i / \sum w_i$

©2005-2007 Carlos Guestrin

17

# Weighting functions

$$w_i = \exp(-D(x_i, \text{query})^2 / K_w^2)$$



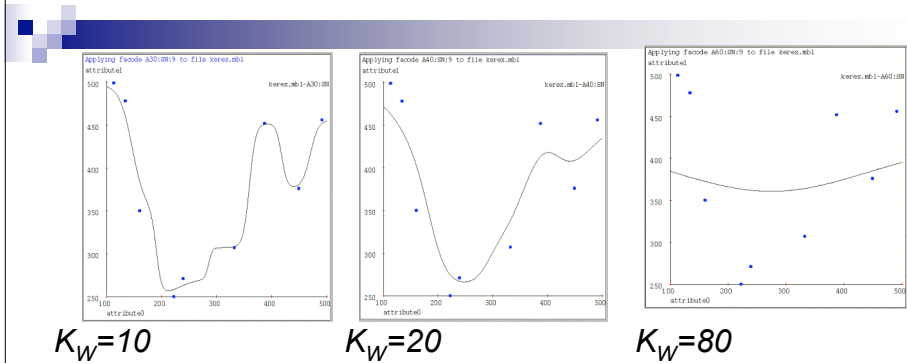
Typically optimize  $K_w$   
using gradient descent

(Our examples use Gaussian)

©2005-2007 Carlos Guestrin

18

# Kernel regression predictions



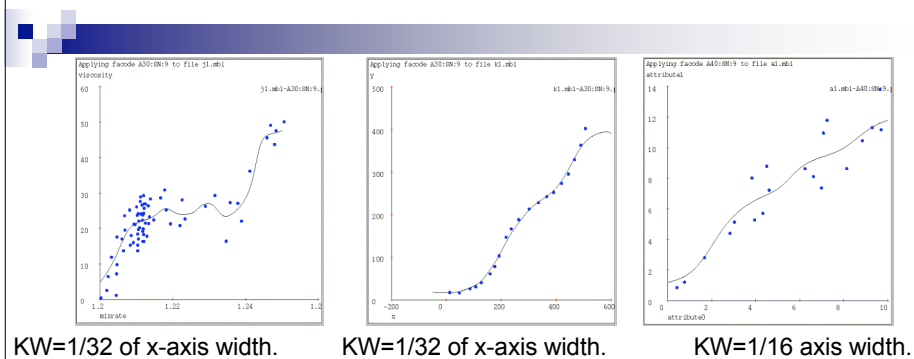
Increasing the kernel width  $K_w$  means further away points get an opportunity to influence you.

As  $K_w \rightarrow \infty$ , the prediction tends to the global average.

©2005-2007 Carlos Guestrin

19

# Kernel regression on our test cases

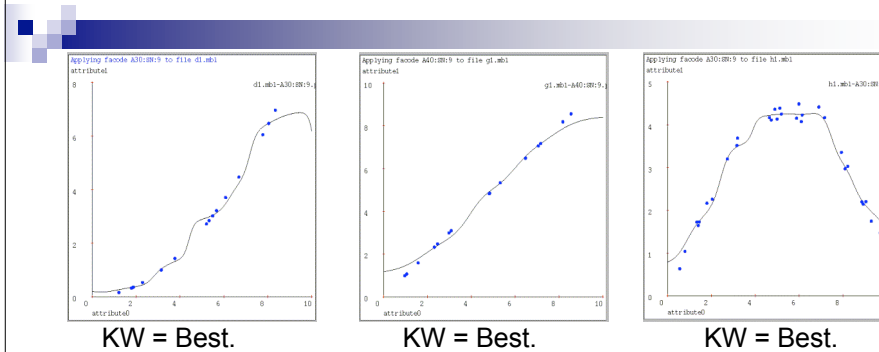


Choosing a good  $K_w$  is important. Not just for Kernel Regression, but for all the locally weighted learners we're about to see.

©2005-2007 Carlos Guestrin

20

## Kernel regression can look bad



Time to try something more powerful...

©2005-2007 Carlos Guestrin

21

## Locally weighted regression

### Kernel regression:

Take a very very conservative function approximator called AVERAGING. Locally weight it.

### Locally weighted regression:

Take a conservative function approximator called LINEAR REGRESSION. Locally weight it.

©2005-2007 Carlos Guestrin

22

# Locally weighted regression

- Four things make a memory based learner:
  - A distance metric
  - Any
  - How many nearby neighbors to look at?
    - All of them
  - A weighting function (optional)
- Kernels
  - $w_i = \exp(-D(x_i, query)^2 / Kw^2)$
- How to fit with the local points?
 

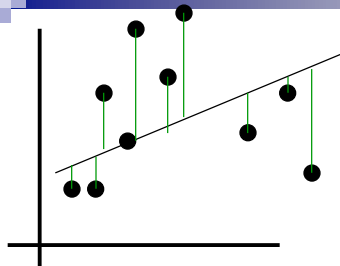
**General weighted regression:**

$$\hat{a} = \underset{\hat{a}}{\operatorname{argmin}} \sum_{k=1}^N w_k^2 (y_k - \hat{a}^T x_k)^2$$

©2005-2007 Carlos Guestrin

23

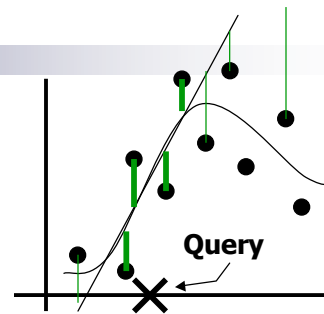
## How LWR works



### Linear regression

- Same parameters for all queries

$$\hat{a} = (X^T X)^{-1} X^T Y$$



### Locally weighted regression

- Solve weighted linear regression for each query

$$\hat{a} = (WX^T WX)^{-1} WX^T WY$$

$$W = \begin{pmatrix} w_1 & 0 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & w_n \end{pmatrix}$$

©2005-2007 Carlos Guestrin

24

## Another view of LWR

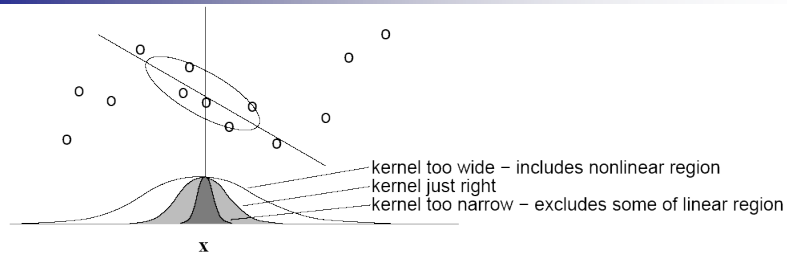
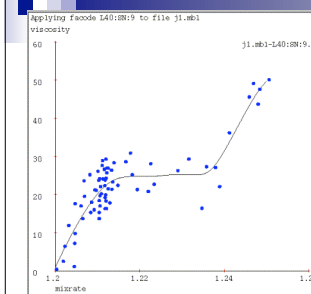
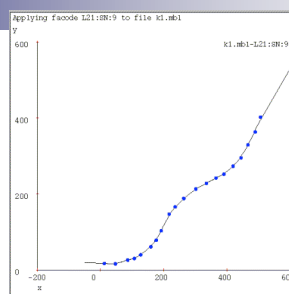


Image from Cohn, D.A., Ghahramani, Z., and Jordan, M. (1996) "Local Learning with Statistical Models", JAIR Volume 4, pages 259-145.

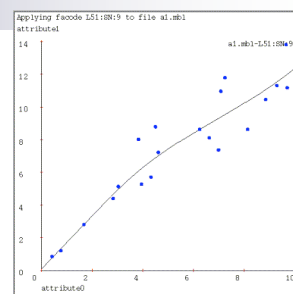
## LWR on our test cases



KW = 1/16 of x-axis width.

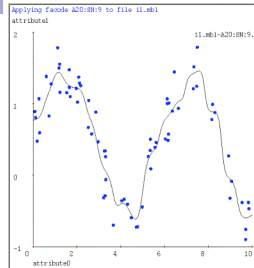


KW = 1/32 of x-axis width.



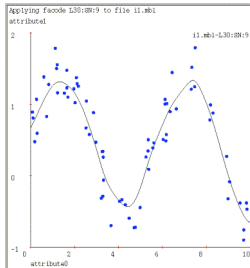
KW = 1/8 of x-axis width.

## Locally weighted polynomial regression



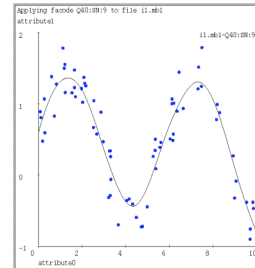
Kernel Regression  
Kernel width  $K_W$  at optimal level.

$KW = 1/100$  x-axis



LW Linear Regression  
Kernel width  $K_W$  at optimal level.

$KW = 1/40$  x-axis



LW Quadratic Regression  
Kernel width  $K_W$  at optimal level.

$KW = 1/15$  x-axis

Local quadratic regression is easy: just add quadratic terms to the WXTWX matrix. As the regression degree increases, the kernel width can increase without introducing bias.

©2005-2007 Carlos Guestrin

27

## Curse of dimensionality for instance-based learning

- Must store and retrieve all data!
  - Most real work done during testing
  - For every test sample, must search through all dataset – very slow!
  - We'll see fast methods for dealing with large datasets
- Instance-based learning often poor with noisy or irrelevant features

©2005-2007 Carlos Guestrin

28

# Curse of the irrelevant feature

©2005-2007 Carlos Guestrin

29

## What you need to know about instance-based learning

- **k-NN**
  - ☐ Simplest learning algorithm
  - ☐ With sufficient data, very hard to beat “strawman” approach
  - ☐ Picking k?
- **Kernel regression**
  - ☐ Set k to n (number of data points) and optimize weights by gradient descent
  - ☐ Smoother than k-NN
- **Locally weighted regression**
  - ☐ Generalizes kernel regression, not just local average
- **Curse of dimensionality**
  - ☐ Must remember (very large) dataset for prediction
  - ☐ Irrelevant features often killers for instance-based approaches

©2005-2007 Carlos Guestrin

30

# Acknowledgment

- This lecture contains some material from Andrew Moore's excellent collection of ML tutorials:

□ <http://www.cs.cmu.edu/~awm/tutorials>

# Support Vector Machines

Machine Learning – 10701/15781

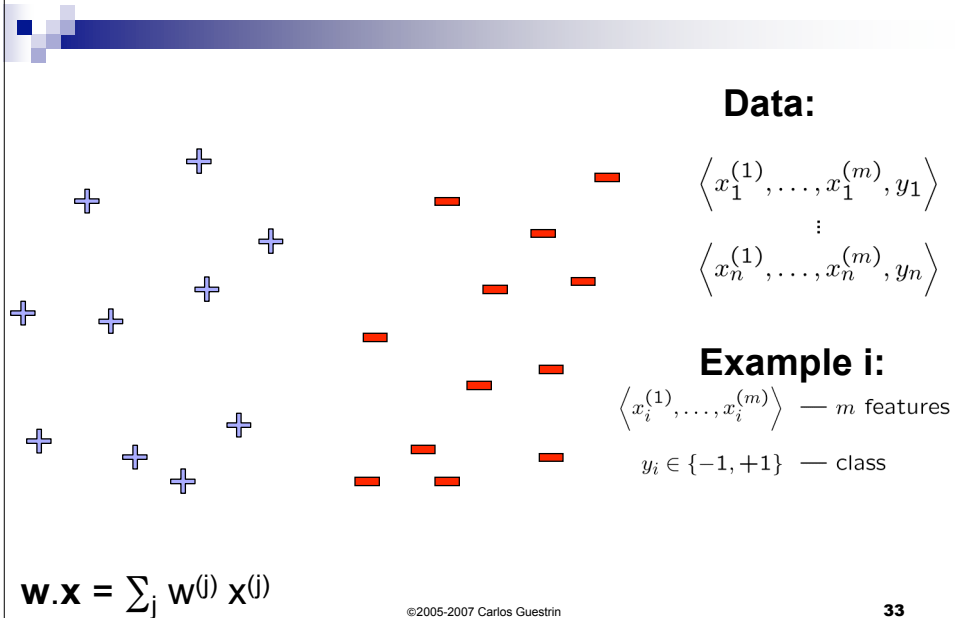
Carlos Guestrin

Carnegie Mellon University

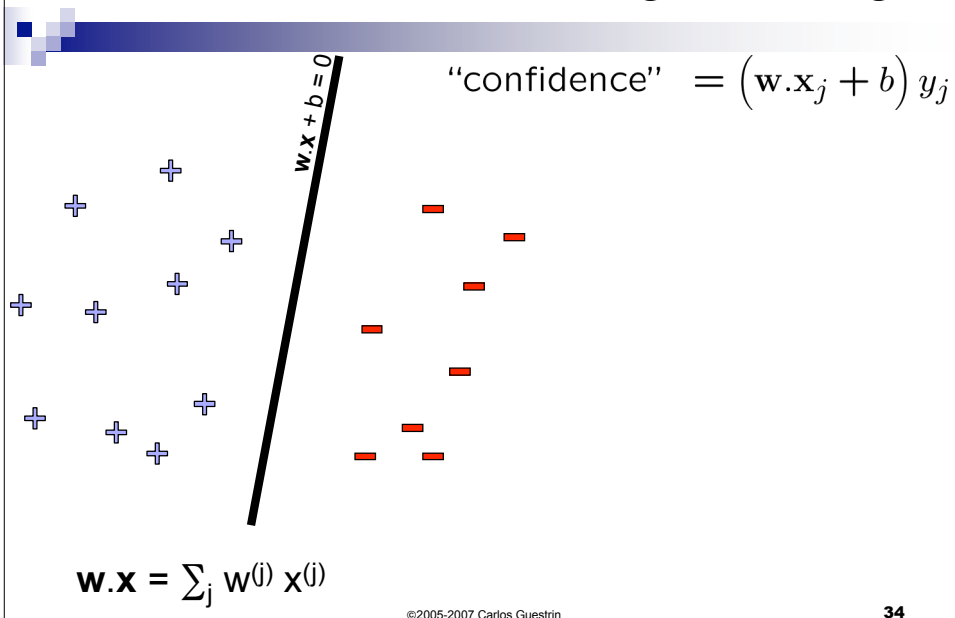
February 19<sup>th</sup>, 2007



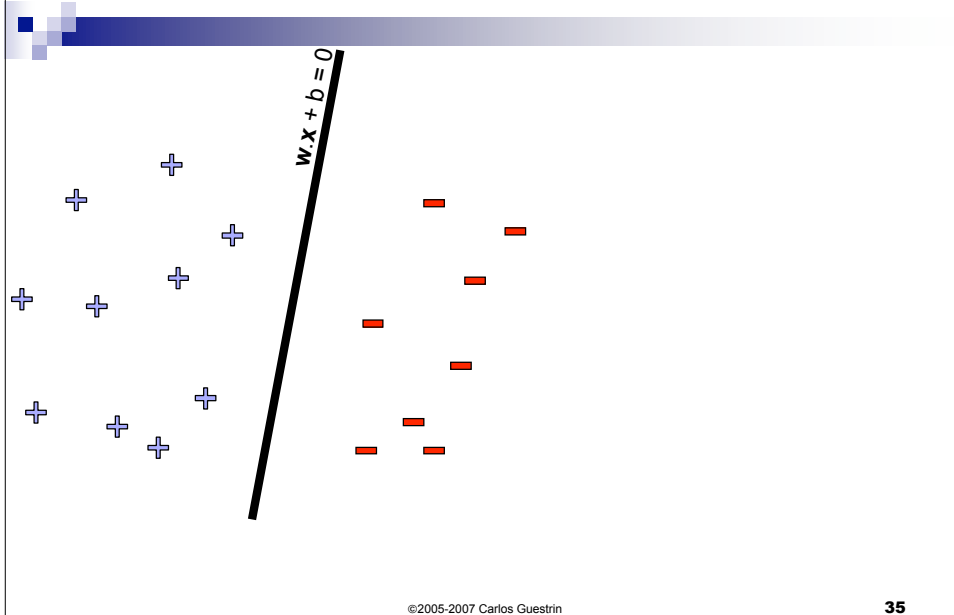
## Linear classifiers – Which line is better?



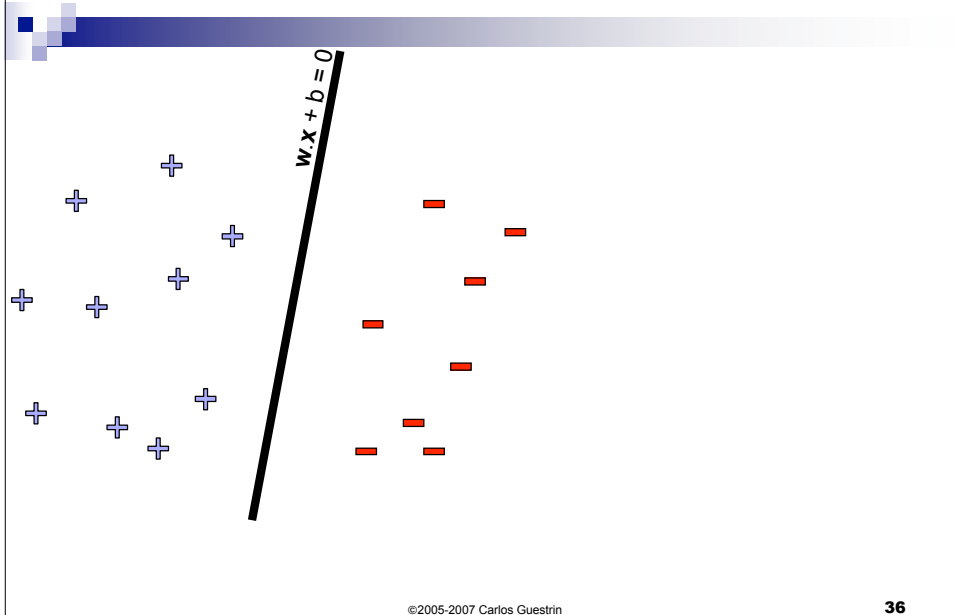
## Pick the one with the largest margin!



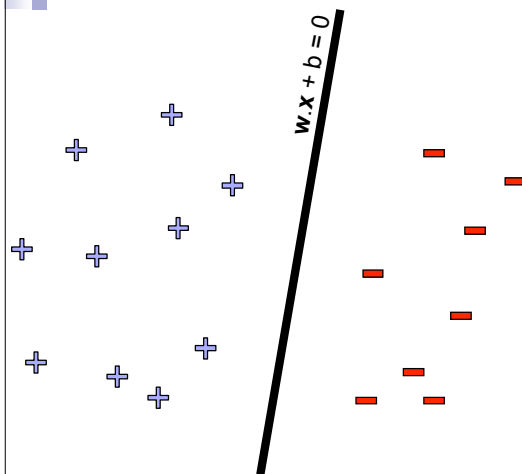
## Maximize the margin



## But there are a many planes...



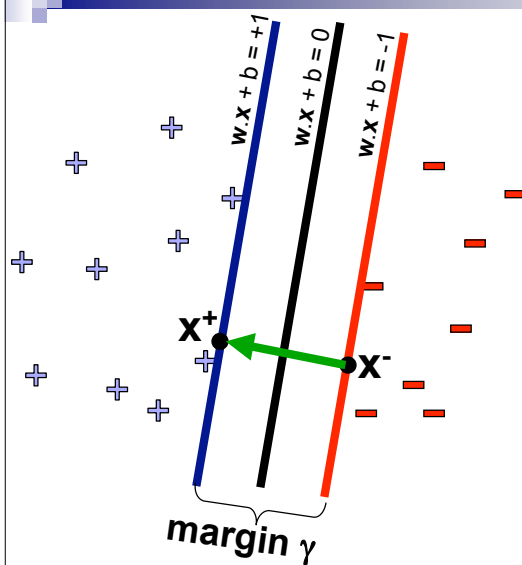
## Review: Normal to a plane



©2005-2007 Carlos Guestrin

37

## Normalized margin – Canonical hyperplanes

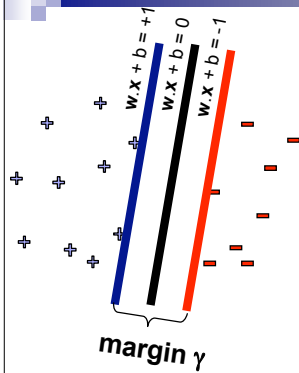


$$\gamma = \frac{2}{\sqrt{w \cdot w}}$$

©2005-2007 Carlos Guestrin

38

## Margin maximization using canonical hyperplanes

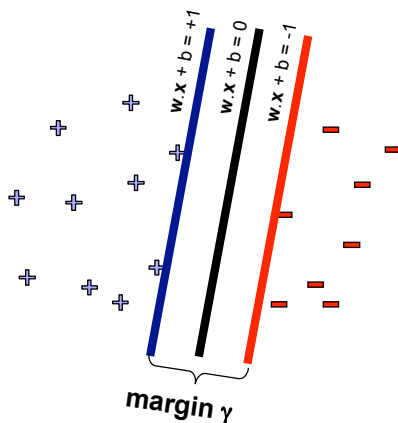


$$\begin{aligned} &\text{minimize}_{\mathbf{w}} \quad \mathbf{w} \cdot \mathbf{w} \\ &(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \quad \forall j \in \text{Dataset} \end{aligned}$$

©2005-2007 Carlos Guestrin

39

## Support vector machines (SVMs)



$$\begin{aligned} &\text{minimize}_{\mathbf{w}} \quad \mathbf{w} \cdot \mathbf{w} \\ &(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \quad \forall j \end{aligned}$$

- Solve efficiently by quadratic programming (QP)
  - Well-studied solution algorithms
- Hyperplane defined by support vectors

©2005-2007 Carlos Guestrin

40