



Instance-based Learning

Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University

October 15th, 2007

©2005-2007 Carlos Guestrin

1

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

2

Consistency of 1-NN

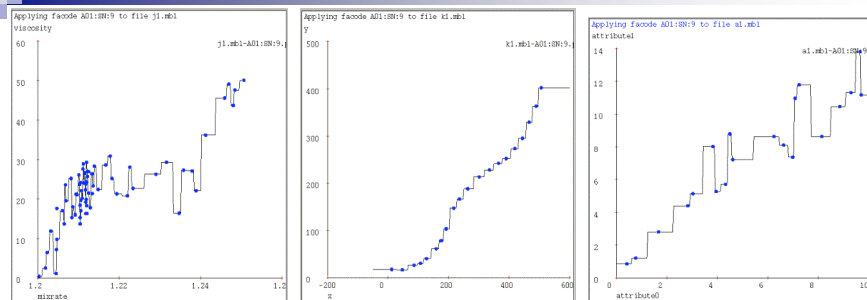
- Consider an estimator f_n trained on n examples
 - e.g., 1-NN, neural nets, regression,...
- Estimator is *consistent* if true error goes to zero as amount of data increases
 - e.g., for no noise data, consistent if:
$$\lim_{n \rightarrow \infty} MSE(f_n) = 0$$
- Regression is not consistent!
 - Representation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

What about variance???

©2005-2007 Carlos Guestrin

3

1-NN overfits?



©2005-2007 Carlos Guestrin

4

k-Nearest Neighbor

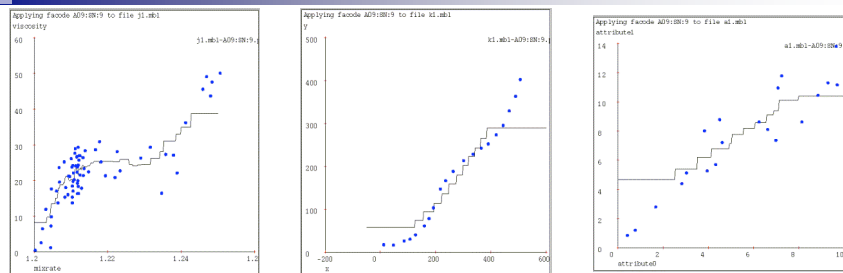
Four things make a memory based learner:

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

5

k-Nearest Neighbor (here k=9)



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

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

©2005-2007 Carlos Guestrin

6

Weighted k-NNs

- Neighbors are not all the same

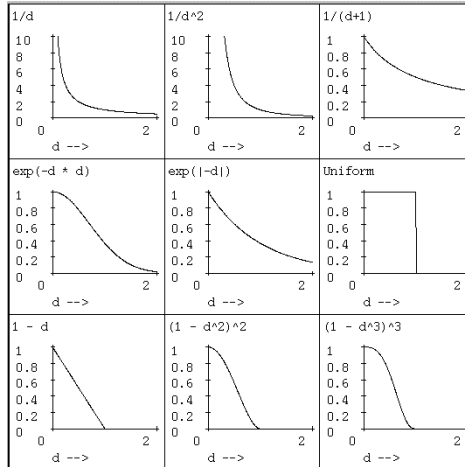
Kernel regression

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)*
 $w_i = \exp(-D(x_i, query)^2 / K_w^2)$
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:
 $predict = \Sigma w_i y_i / \Sigma w_i$

Weighting functions

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



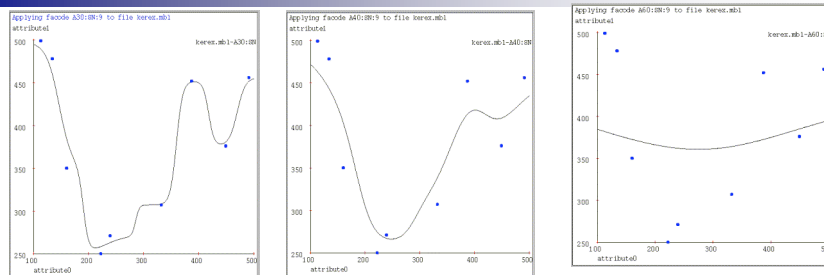
Typically optimize K_w using gradient descent

(Our examples use Gaussian)

©2005-2007 Carlos Guestrin

9

Kernel regression predictions



$K_w=10$

$K_w=20$

$K_w=80$

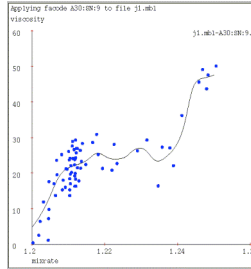
Increasing the kernel width K_w means further away points get an opportunity to influence you.

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

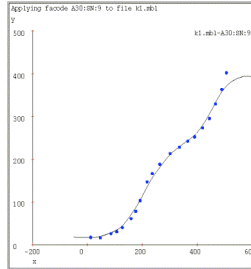
©2005-2007 Carlos Guestrin

10

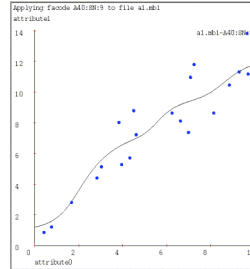
Kernel regression on our test cases



KW=1/32 of x-axis width.



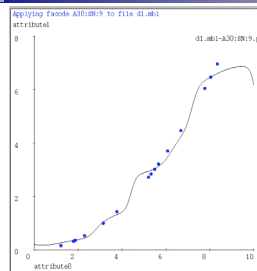
KW=1/32 of x-axis width.



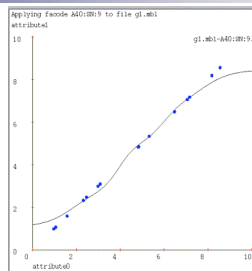
KW=1/16 axis width.

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

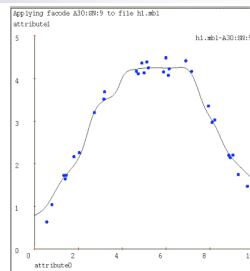
Kernel regression can look bad



KW = Best.



KW = Best.



KW = Best.

Time to try something more powerful...

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.

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

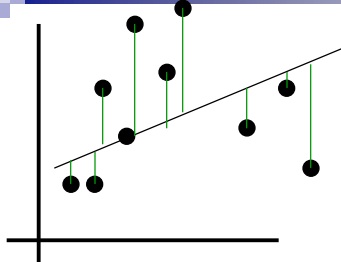
- $w_i = \exp(-D(x_i, query)^2 / Kw^2)$

- *How to fit with the local points?*

General weighted regression:

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

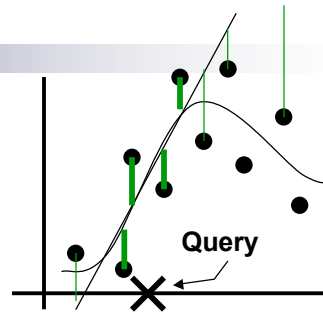
How LWR works



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

15

Another view of LWR

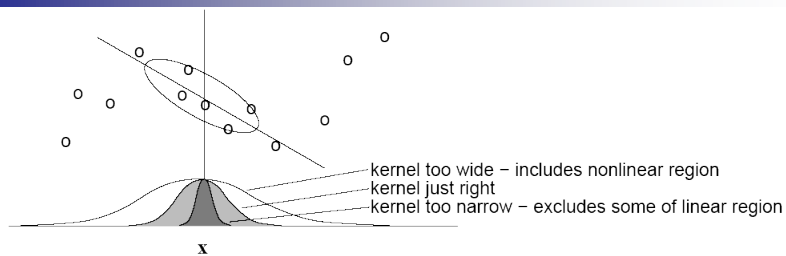
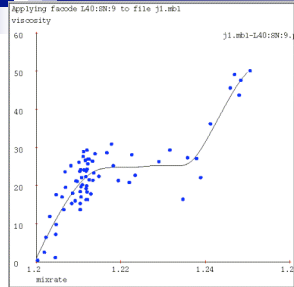
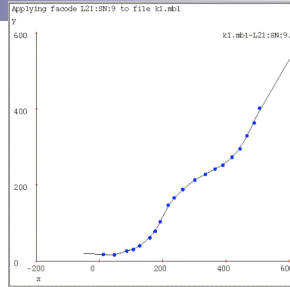


Image from Cohn, D.A., Ghahramani, Z., and Jordan, M.I. (1996). "Learning with Statistical Models", JAIR Volume 4, pages 109-145.

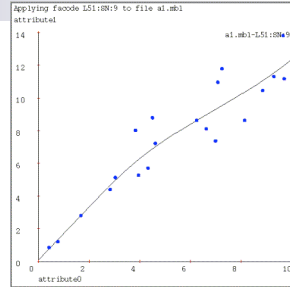
LWR on our test cases



KW = 1/16 of x-axis width.



KW = 1/32 of x-axis width.

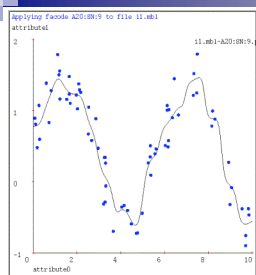


KW = 1/8 of x-axis width.

©2005-2007 Carlos Guestrin

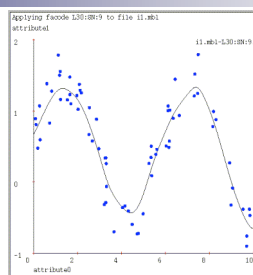
17

Locally weighted polynomial regression



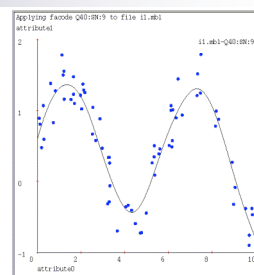
Kernel Regression
Kernel width K_W at optimal level.

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

18

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

Curse of the irrelevant feature

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

Acknowledgment

- This lecture contains some material from Andrew Moore’s excellent collection of ML tutorials:
 - <http://www.cs.cmu.edu/~awm/tutorials>

Announcements

- Recitation this week: Neural networks

- Project proposals due next Wednesday
 - Exciting data:
 - Swivel.com - user generated graphs
 - Recognizing Captchas
 - Election contributions
 - Activity recognition
 - ...

Support Vector Machines

Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University

October 15th, 2007

Linear classifiers – Which line is better?

Data:

$$\langle x_1^{(1)}, \dots, x_1^{(m)}, y_1 \rangle$$

$$\vdots$$

$$\langle x_n^{(1)}, \dots, x_n^{(m)}, y_n \rangle$$

Example i:

$$\langle x_i^{(1)}, \dots, x_i^{(m)} \rangle \text{ — } m \text{ features}$$

$$y_i \in \{-1, +1\} \text{ — class}$$

$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$

©2005-2007 Carlos Guestrin 25

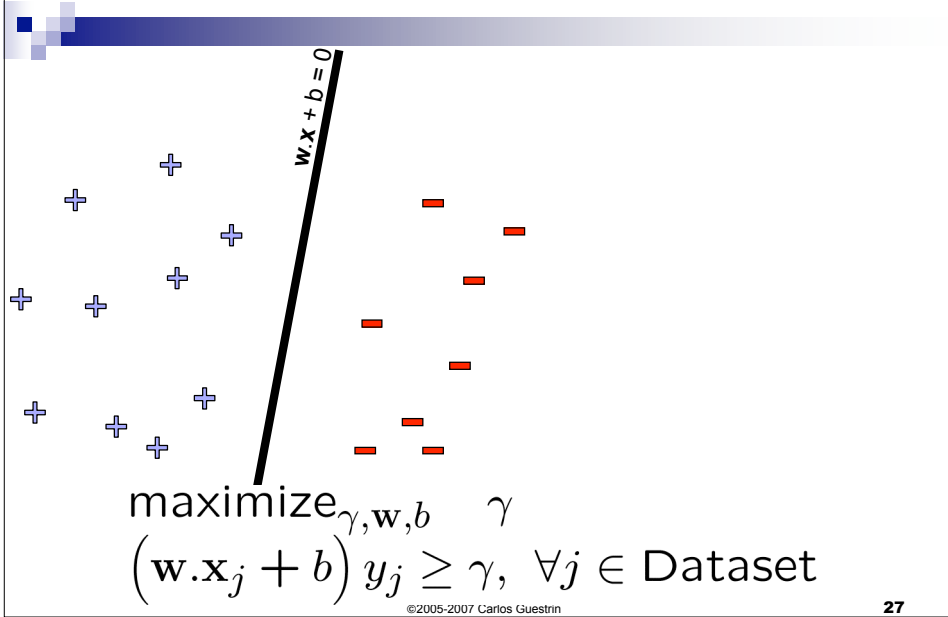
Pick the one with the largest margin!

“confidence” = $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j$

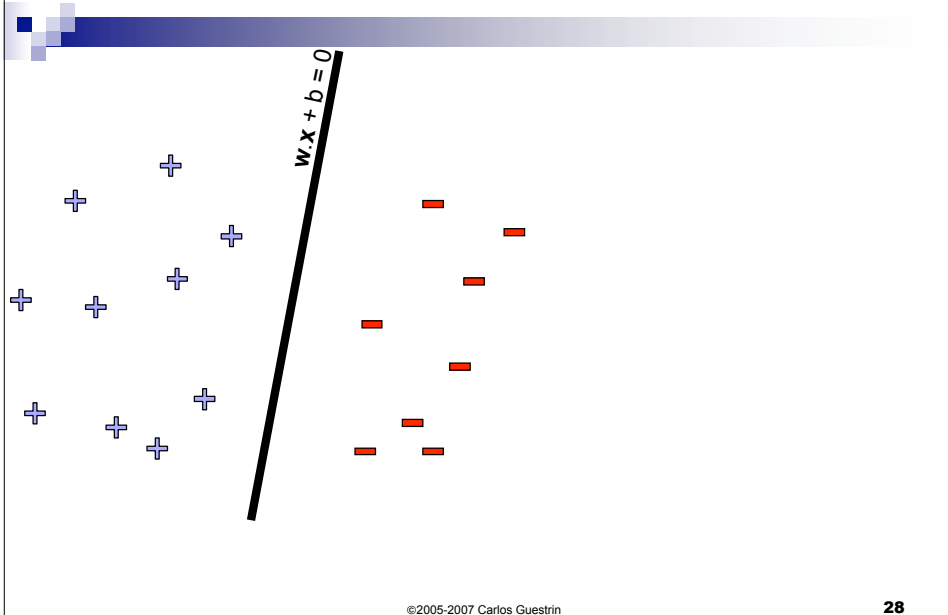
$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$

©2005-2007 Carlos Guestrin 26

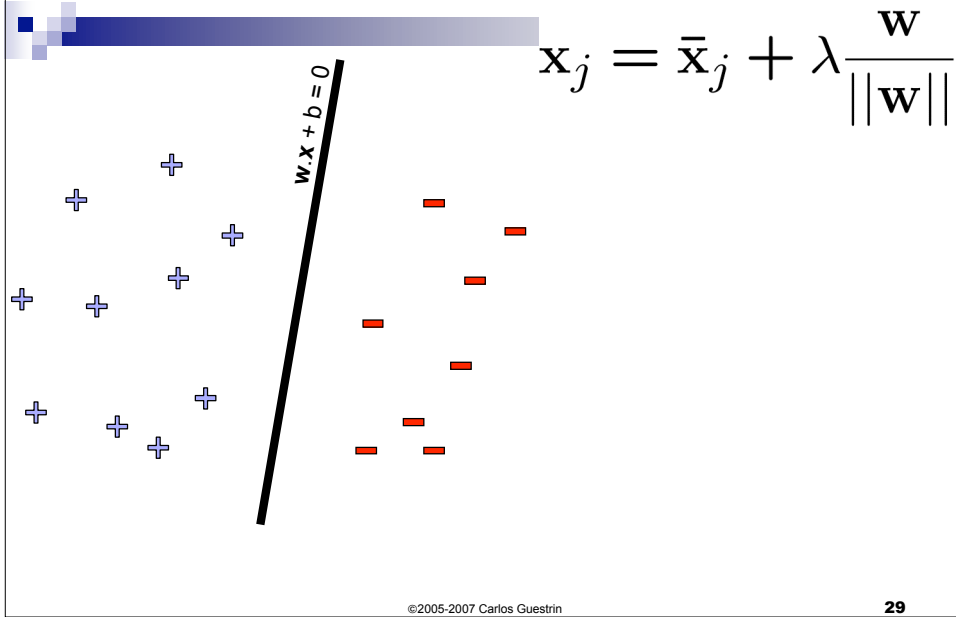
Maximize the margin



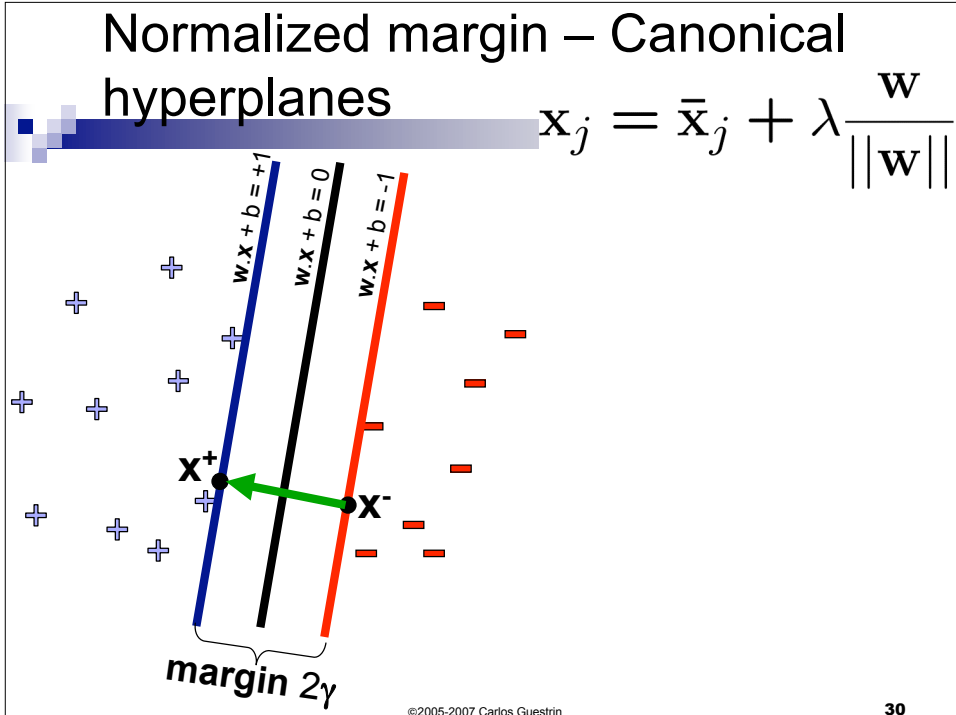
But there are a many planes...



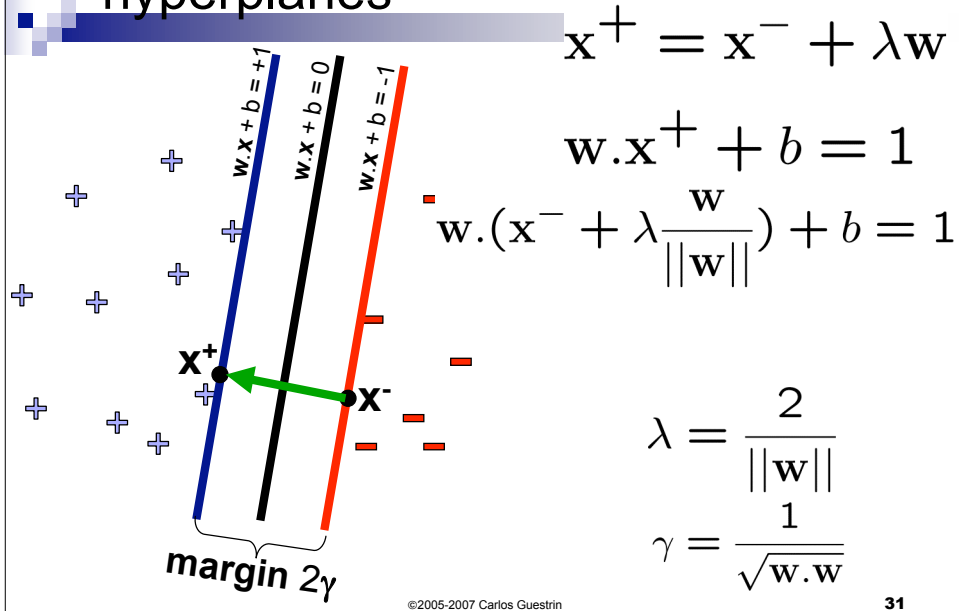
Review: Normal to a plane



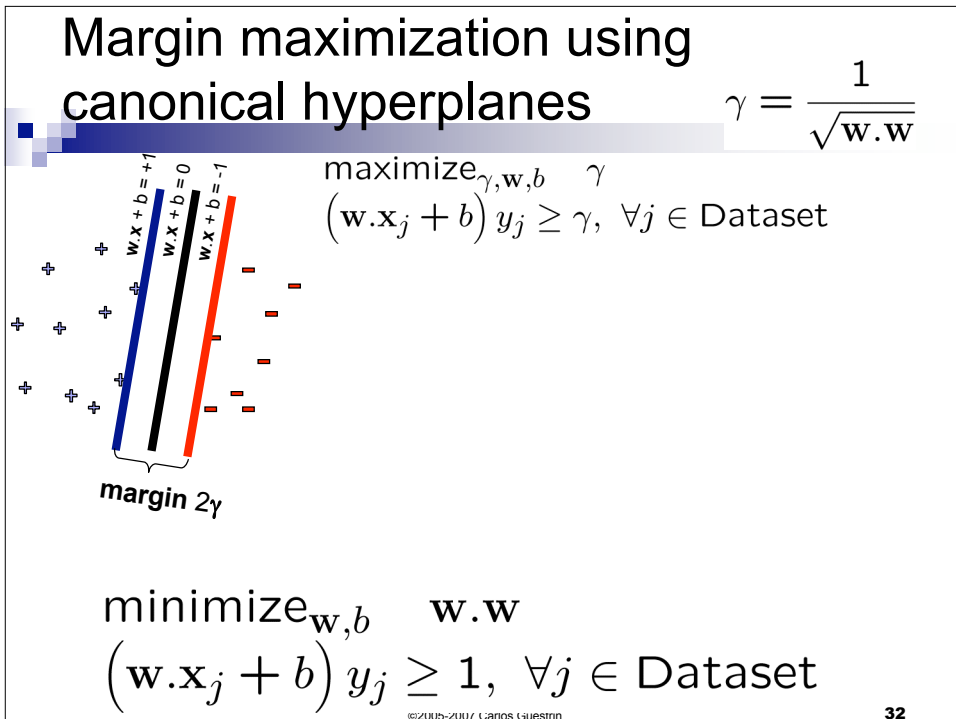
Normalized margin – Canonical hyperplanes



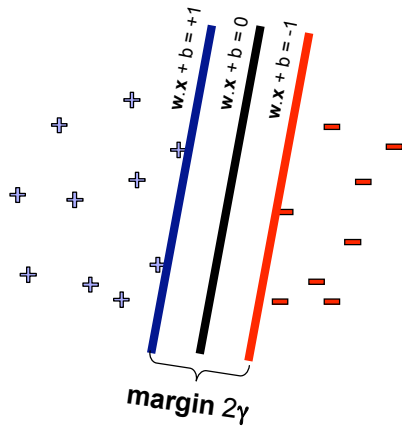
Normalized margin – Canonical hyperplanes



Margin maximization using canonical hyperplanes



Support vector machines (SVMs)



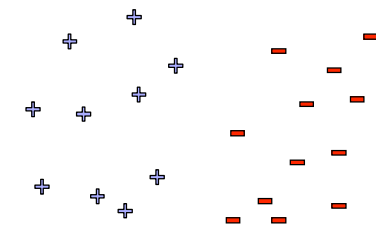
$$\text{minimize}_{w,b} \quad w \cdot w$$
$$\left(w \cdot x_j + b \right) y_j \geq 1, \quad \forall j$$

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

©2005-2007 Carlos Guestrin

33

What if the data is not linearly separable?



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

©2005-2007 Carlos Guestrin

34

What if the data is still not linearly separable?

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

- Minimize $\mathbf{w} \cdot \mathbf{w}$ and number of training mistakes
 - Tradeoff two criteria?
- Tradeoff #(mistakes) and $\mathbf{w} \cdot \mathbf{w}$
 - 0/1 loss
 - Slack penalty C
 - Not QP anymore
 - Also doesn't distinguish near misses and really bad mistakes

©2005-2007 Carlos Guestrin 35

Slack variables – Hinge loss

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

- If margin ≥ 1 , don't care
- If margin < 1 , pay linear penalty

©2005-2007 Carlos Guestrin 36

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

SVM:

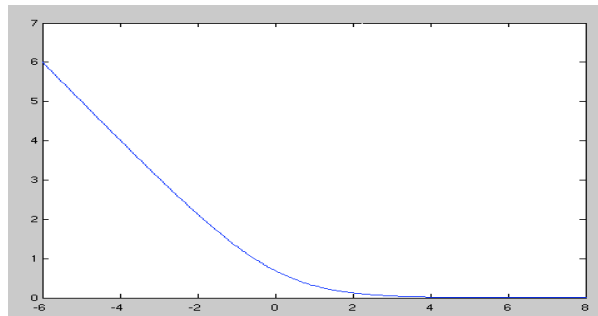
$$\begin{aligned} &\text{minimize}_{w,b} \quad w \cdot w + C \sum_j \xi_j \\ &(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j \\ &\xi_j \geq 0, \quad \forall j \end{aligned}$$

Logistic regression:

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

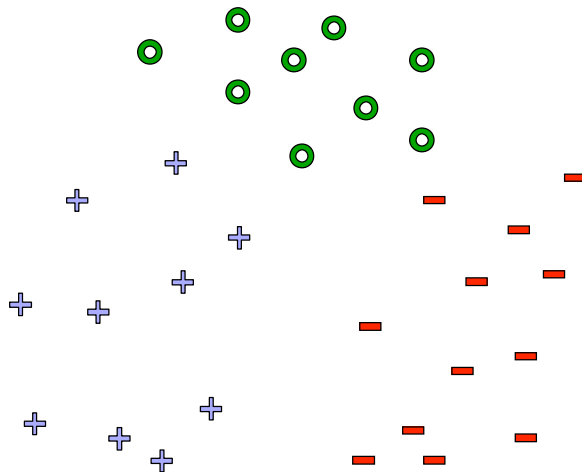
Log loss:

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



37

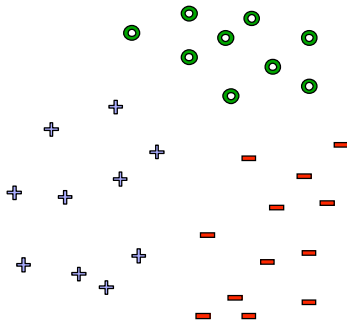
What about multiple classes?



38

One against All

Learn 3 classifiers:

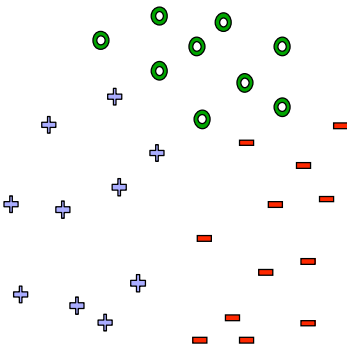


©2005-2007 Carlos Guestrin

39

Learn 1 classifier: Multiclass SVM

Simultaneously learn 3 sets of weights



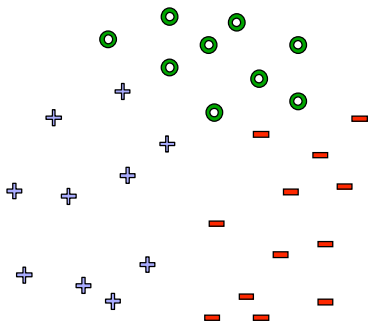
$$w^{(y_j)} \cdot x_j + b^{(y_j)} \geq w^{(y')} \cdot x_j + b^{(y')} + 1, \quad \forall y' \neq y_j, \quad \forall j$$

©2005-2007 Carlos Guestrin

40

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

41

What you need to know

- Maximizing margin
- Derivation of SVM formulation
- Slack variables and hinge loss
- Relationship between SVMs and logistic regression
 - 0/1 loss
 - Hinge loss
 - Log loss
- Tackling multiple class
 - One against All
 - Multiclass SVMs

©2005-2007 Carlos Guestrin

42