

Neural Networks, *cont.*

Machine Learning – 10701/15781

Carlos Guestrin

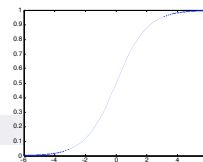
Carnegie Mellon University

October 10th, 2007

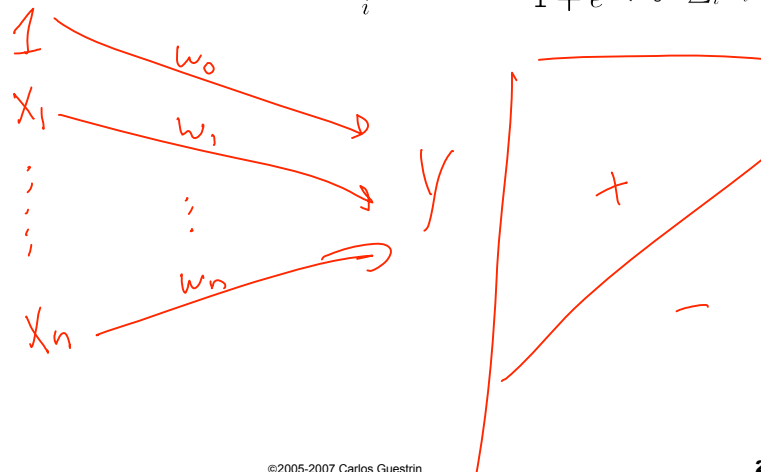
©2005-2007 Carlos Guestrin

1

Perceptron as a graph



$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$



©2005-2007 Carlos Guestrin

2

The perceptron learning rule

Squared loss

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j \delta^j$$

learning rate

$$\delta^j = [y^j - g(w_0 + \sum_i w_i x_i^j)] g^j (1 - g^j)$$

$$g^j = g(w_0 + \sum_i w_i x_i^j)$$

■ Compare to MLE:

$$w_i \leftarrow w_i + \eta \sum_j x_i^j \delta^j \quad \delta^j = [y^j - g(w_0 + \sum_i w_i x_i^j)]$$

©2005-2007 Carlos Guestrin 3

Hidden layer

■ Perceptron: $out(x) = g(w_0 + \sum_i w_i x_i)$

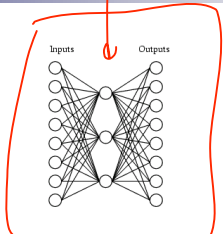
■ 1-hidden layer:

$$out(x) = g\left(w_0 + \sum_k w_k g\left(w_0^k + \sum_i w_i^k x_i\right)\right)$$

Diagram illustrating a 1-hidden layer perceptron:

©2005-2007 Carlos Guestrin 4

Example data for NN with hidden layer



A target function: $x \rightarrow \text{out}(x)$

dataset

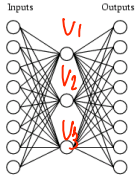
Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

Can this be learned??

5

Learned weights for hidden layer

A network:



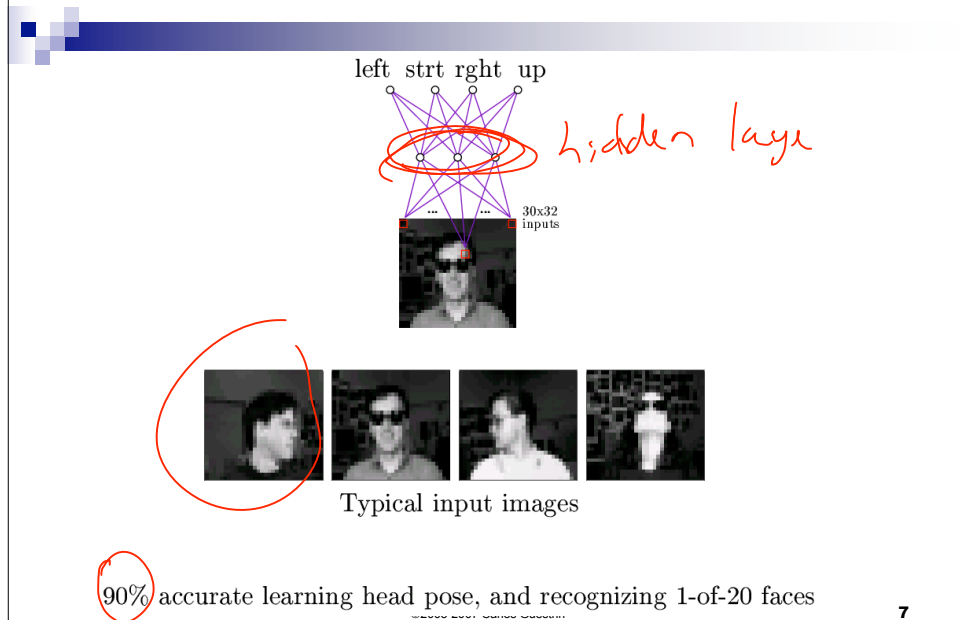
Learned hidden layer representation:

Input	Hidden Values	Output
10000000	→ .89 .04 .08	→ 10000000
01000000	→ .01 .11 .88	→ 01000000
00100000	→ .01 .97 .27	→ 00100000
00010000	→ .99 .97 .71	→ 00010000
00001000	→ .03 .05 .02	→ 00001000
00000100	→ .22 .99 .99	→ 00000100
00000010	→ .80 .01 .98	→ 00000010
00000001	→ .60 .94 .01	→ 00000001

"learned" a lower dimensional encoding of the data

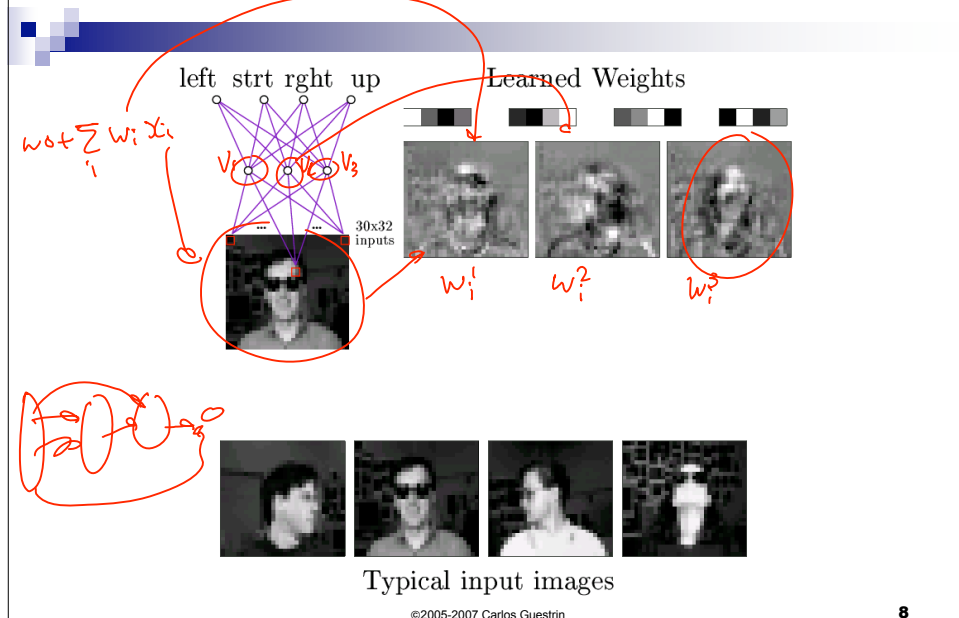
6

NN for images



7

Weights in NN for images



8

Gradient descent for 1-hidden layer – Back-propagation: Computing $\frac{\partial \ell(W)}{\partial w_k}$

loss
 $\ell(W) = \frac{1}{2} \sum_j [y^j - \text{out}(\mathbf{x}^j)]^2$ *doesn't depend on w_k*

data point
 $\text{out}(\mathbf{x}) = g\left(\sum_{k'} w_{k'} g\left(\sum_{i'} w_{i'}^{k'} x_{i'}\right)\right)$

Dropped w_0 to make derivation simpler

$\frac{\partial \ell(W)}{\partial w_k} = \sum_{j=1}^m -[y^j - \text{out}(\mathbf{x}^j)] \frac{\partial \text{out}(\mathbf{x}^j)}{\partial w_k}$ *$\frac{\partial g(f(w))}{\partial w} = \frac{\partial f}{\partial w} \cdot g'(f(w))$*

$\frac{\partial \text{out}(\mathbf{x}^j)}{\partial w_k} = V_k^j \cdot g'\left(\sum_{k'} w_{k'} g\left(\sum_{i'} w_{i'}^{k'} x_{i'}^j\right)\right)$ *$g' = g(1-g)$*

$= V_k^j g\left(\sum_{k'} w_{k'} V_{k'}^j\right) (1 - g\left(\sum_{k'} w_{k'} V_{k'}^j\right))$

first compute V_k 's (forward sweep)

©2005-2007 Carlos Guestrin

9

Gradient descent for 1-hidden layer – Back-propagation: Computing $\frac{\partial \ell(W)}{\partial w_i^k}$

Dropped w_0 to make derivation simpler

$\ell(W) = \frac{1}{2} \sum_j [y^j - \text{out}(\mathbf{x}^j)]^2$

$\text{out}(\mathbf{x}) = g\left(\sum_{k'} w_{k'} g\left(\sum_{i'} w_{i'}^{k'} x_{i'}\right)\right)$

$\frac{\partial \ell(W)}{\partial w_i^k} = \sum_{j=1}^m -[y^j - \text{out}(\mathbf{x}^j)] \frac{\partial \text{out}(\mathbf{x}^j)}{\partial w_i^k}$ *$\frac{\partial g(f(h(w)))}{\partial w} = \frac{\partial g}{\partial f} \cdot \frac{\partial f}{\partial h} \cdot \frac{\partial h}{\partial w}$*

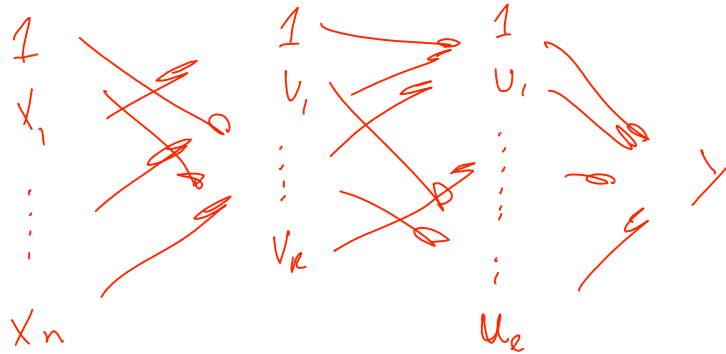
$\frac{\partial \text{out}(\mathbf{x}^j)}{\partial w_i^k} = x_i^j w_k g'\left(\sum_{i'} w_{i'}^k x_{i'}^j\right) \cdot g'\left(\sum_{k'} w_{k'} g\left(\sum_{i'} w_{i'}^{k'} x_{i'}^j\right)\right)$

$= x_i^j w_k V_k^j (1 - V_k^j) \cdot \frac{g\left(\sum_{k'} w_{k'} V_{k'}^j\right)}{\text{out}(\mathbf{x}^j)} \cdot \frac{(1 - g\left(\sum_{k'} w_{k'} V_{k'}^j\right))}{(1 - \text{out}(\mathbf{x}^j))}$

©2005-2007 Carlos Guestrin

10

Multilayer neural networks



©2005-2007 Carlos Guestrin

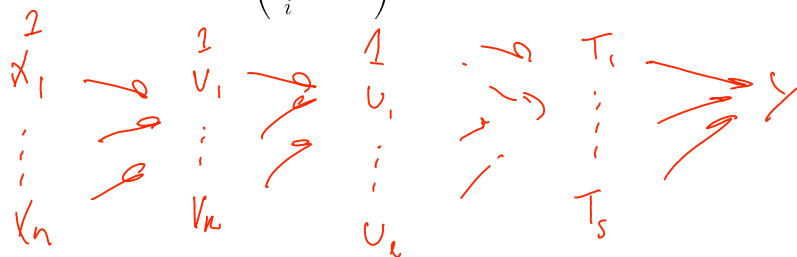
11

Forward propagation – prediction

- Recursive algorithm
- Start from input layer
- Output of node V_k with parents U_1, U_2, \dots :

recurrent NNS (cycles in graph)

$$V_k = g\left(\sum_i w_i^k U_i\right)$$



©2005-2007 Carlos Guestrin

12

Back-propagation – learning

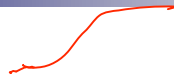
- Just gradient descent!!!
- Recursive algorithm for computing gradient
- For each example
 - Perform forward propagation
 - Start from output layer
 - Compute gradient of node V_k with parents U_1, U_2, \dots
 - Update weight w_i^k

©2005-2007 Carlos Guestrin

13

Many possible response functions

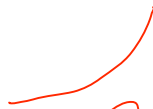
- Sigmoid



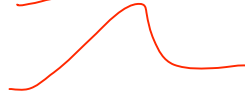
- Linear

$$g(\sum_i w_i x_i) = w_0 + \sum_i w_i x_i$$

- Exponential



- Gaussian



- ...

step/threshold :



$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

©2005-2007 Carlos Guestrin

14

Convergence of backprop

■ Perceptron leads to convex optimization

- Gradient descent reaches **global minima**

simulated annealing, random restarts, prayer, ...

■ Multilayer neural nets **not convex**

- Gradient descent gets stuck in local minima
- Hard to set learning rate
- Selecting number of hidden units and layers = fuzzy process
- NNs falling in disfavor in last few years
- We'll see later in semester, kernel trick is a good alternative
- Nonetheless, neural nets are one of the most used ML approaches

©2005-2007 Carlos Guestrin

15

Overfitting?

■ Neural nets represent complex functions

- Output becomes more complex with gradient steps



©2005-2007 Carlos Guestrin

16

Overfitting

- Output fits training data “too well”
 - Poor test set accuracy
- Overfitting the training data
 - Related to bias-variance tradeoff
 - One of central problems of ML
- Avoiding overfitting?
 - More training data
 - Regularization
 - Early stopping

©2005-2007 Carlos Guestrin

17

What you need to know about neural networks

- Perceptron:
 - Representation
 - Perceptron learning rule
 - Derivation
- Multilayer neural nets
 - Representation
 - Derivation of backprop
 - Learning rule
- Overfitting
 - Definition
 - Training set versus test set
 - Learning curve

©2005-2007 Carlos Guestrin

18

Announcements

- Recitation this week: Neural networks

- Project proposals due next Wednesday

- Exciting data:

- Swivel.com - user generated graphs
- Recognizing Captchas
- Election contributions
- Activity recognition
- ...

*ML is real data
the more, the better*

Instance-based Learning

(non-parametric methods)

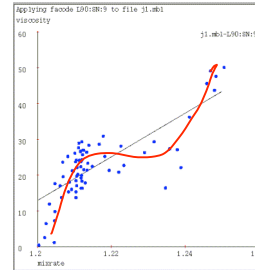
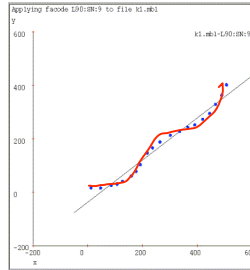
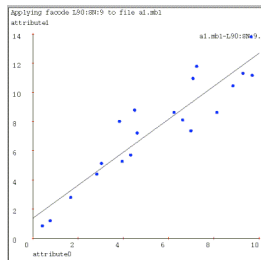
Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

October 10th, 2007

Why not just use Linear Regression?

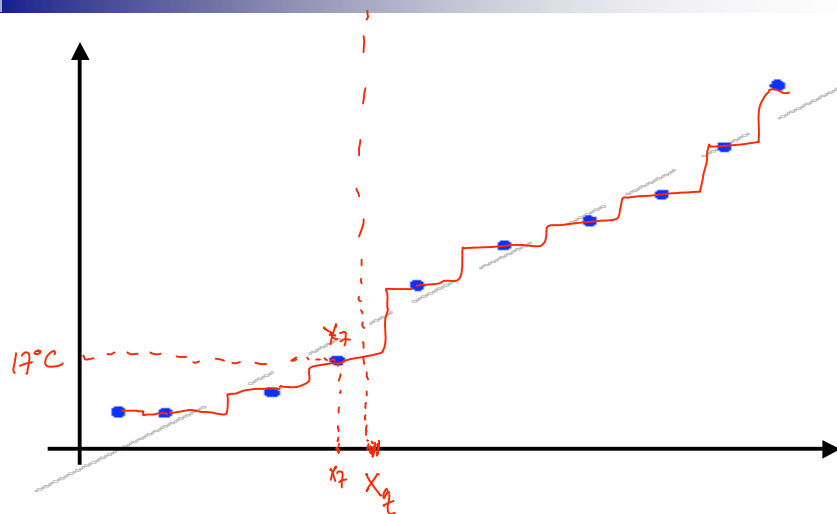


- use more basis functions
- non-parametric ~

©2005-2007 Carlos Guestrin

21

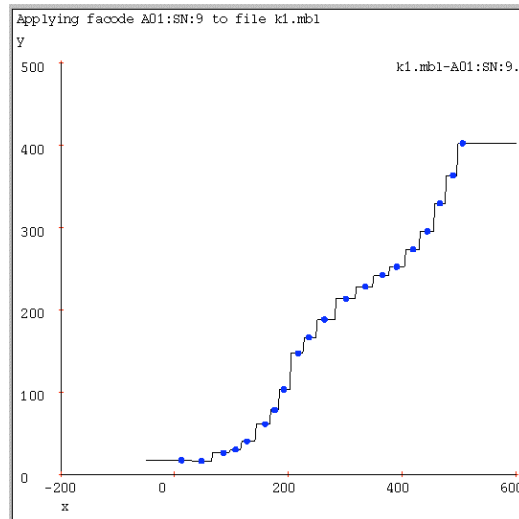
Using data to predict new data



©2005-2007 Carlos Guestrin

22

Nearest neighbor



works amazingly
well with
a lot of
data

©2005-2007 Carlos Guestrin

23

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 $\hat{y} \approx f(x_q)$

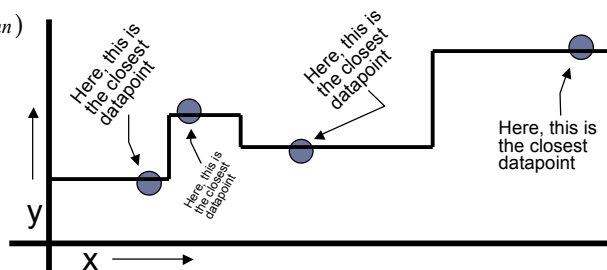
Nearest Neighbor:

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.



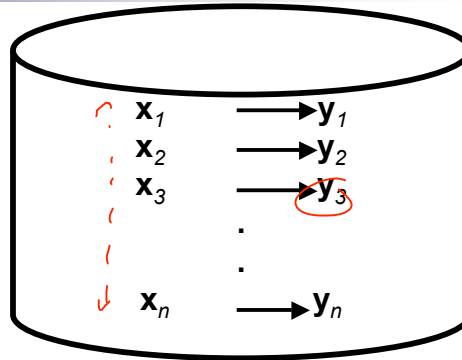
©2005-2007 Carlos Guestrin

24

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 (what's near?)
- How many nearby neighbors to look at?
- A weighting function (optional)
- How to fit with the local points?

©2005-2007 Carlos Guestrin

25

1-Nearest Neighbor

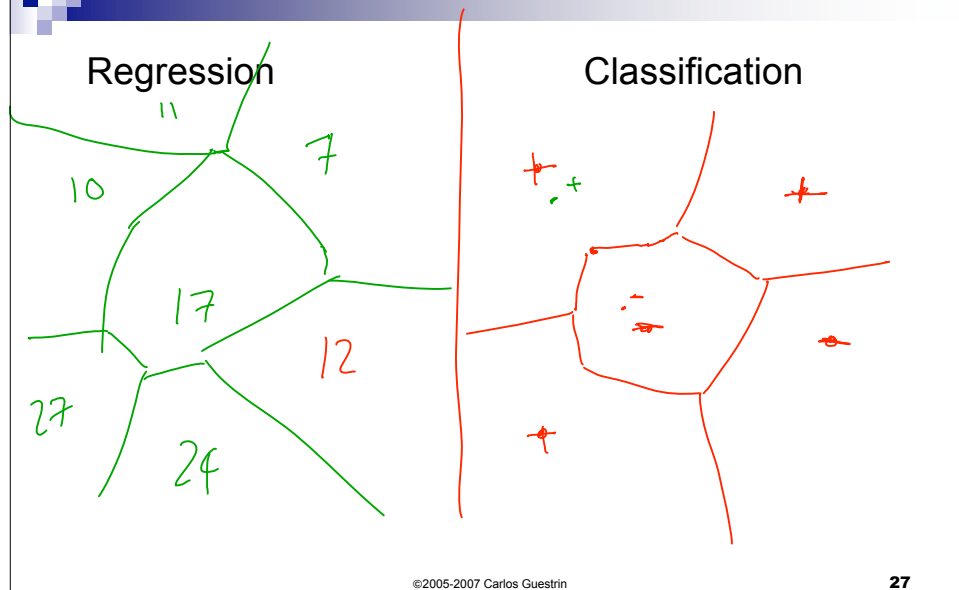
Four things make a memory based learner:

1. A distance metric
Euclidian (and many more) $\|x_q - x_i\|_2$
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

26

Multivariate 1-NN examples

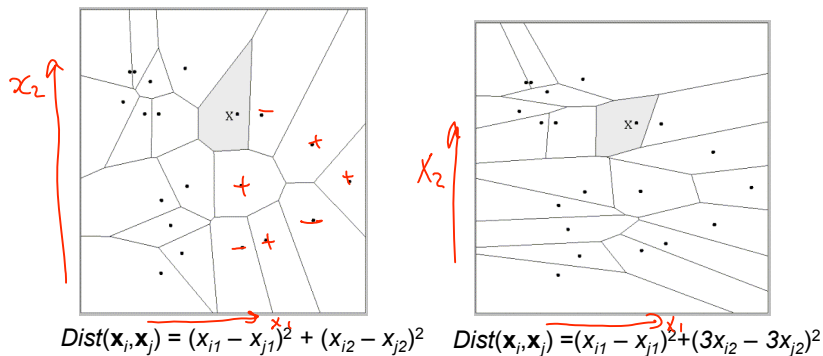


Multivariate distance metrics

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.



The relative scalings in the distance metric affect region shapes

Euclidean distance metric

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

where

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

$\sigma_1 = \sigma_2$ $\sigma_1 > \sigma_2$

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

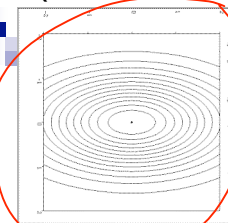
Other Metrics...

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

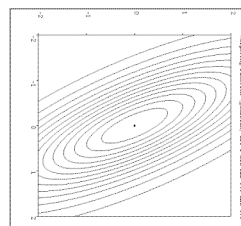
©2005-2007 Carlos Guestrin

29

Notable distance metrics (and their level sets)

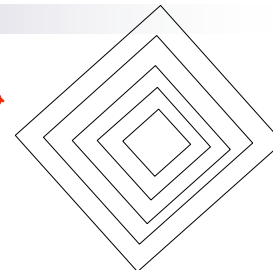


Scaled Euclidean (L_2)



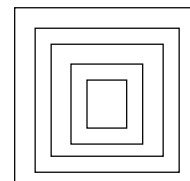
Mahalanobis (here, Σ on the previous slide is not necessarily diagonal, but is symmetric)

$$\|x_1 - x_2\|_1 = \sum_i |x_{i1} - x_{i2}|$$



L_1 norm (absolute)

$$\|x_1 - x_2\|_\infty = \max_i |x_{i1} - x_{i2}|$$



L_∞ (max) norm

©2005-2007 Carlos Guestrin

30

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 (i.i.d.)

- e.g., for no noise data, consistent if:

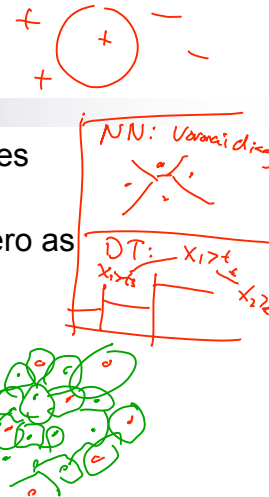
$$\lim_{n \rightarrow \infty} MSE(f_n) = 0$$

- Regression is not consistent!

- Representation bias

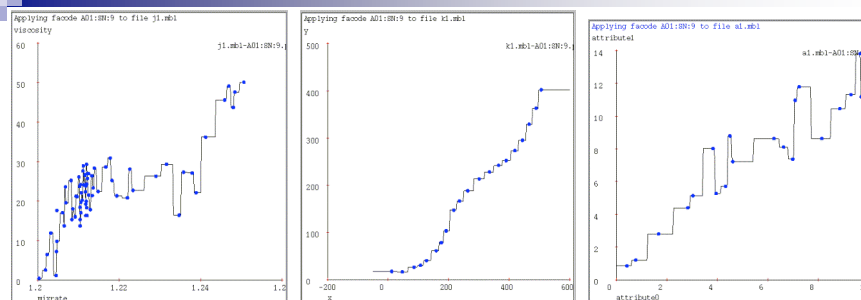
- 1-NN is consistent** (under some mild fineprint)

What about variance???



31

1-NN overfits?



©2005-2007 Carlos Guestrin

32

k-Nearest Neighbor

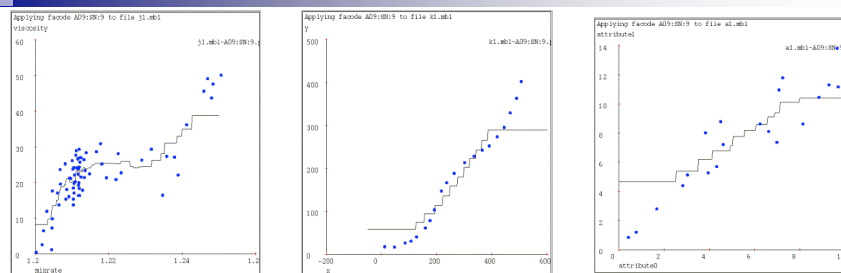
Four things make a memory based learner:

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

33

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

34

Weighted k-NNs

- Neighbors are not all the same

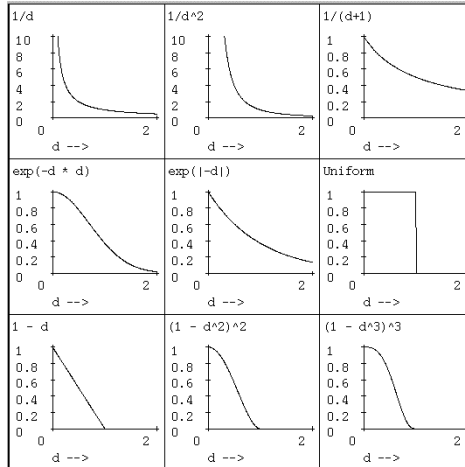
Kernel regression

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:
 $\text{predict} = \sum w_i y_i / \sum w_i$

Weighting functions

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



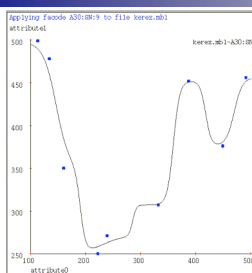
Typically optimize K_w
using gradient descent

(Our examples use Gaussian)

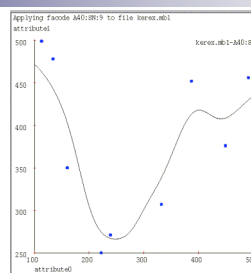
©2005-2007 Carlos Guestrin

37

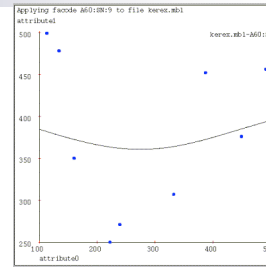
Kernel regression predictions



$K_w=10$



$K_w=20$



$K_w=80$

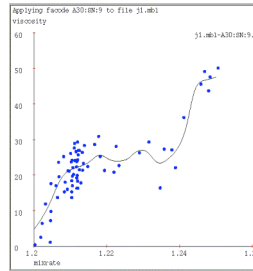
Increasing the kernel width K_w means further away points get an opportunity to influence you.

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

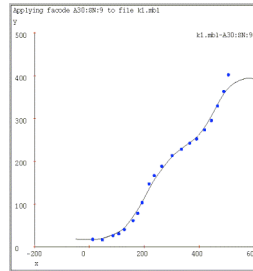
©2005-2007 Carlos Guestrin

38

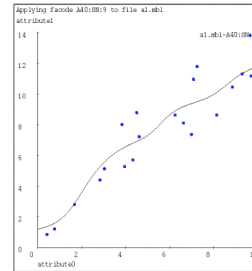
Kernel regression on our test cases



KW=1/32 of x-axis width.



KW=1/32 of x-axis width.



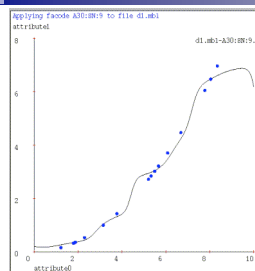
KW=1/16 axis width.

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

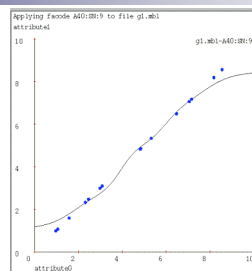
©2005-2007 Carlos Guestrin

39

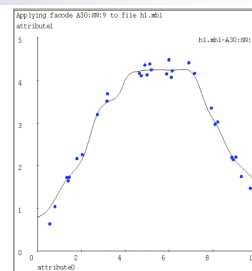
Kernel regression can look bad



KW = Best.



KW = Best.



KW = Best.

Time to try something more powerful...

©2005-2007 Carlos Guestrin

40

Locally weighted regression

Kernel regression:

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

Locally weighted regression:

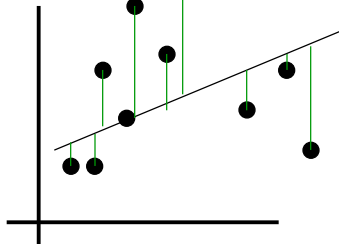
Take a conservative function approximator called LINEAR REGRESSION. Locally weight it.

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
 - $w_i = \exp(-D(x_i, \text{query})^2 / K w^2)$
- *How to fit with the local points?*
General weighted regression:

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

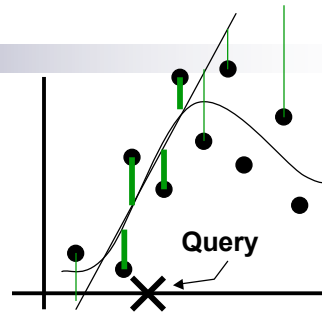
How LWR works



Linear regression

- Same parameters for all queries

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



Locally weighted regression

- Solve weighted linear regression for each query

$$\beta = ((\mathbf{W}\mathbf{X})^T \mathbf{W}\mathbf{X})^{-1} (\mathbf{W}\mathbf{X})^T \mathbf{W}\mathbf{Y}$$

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

©2005-2007 Carlos Guestrin

43

Another view of LWR

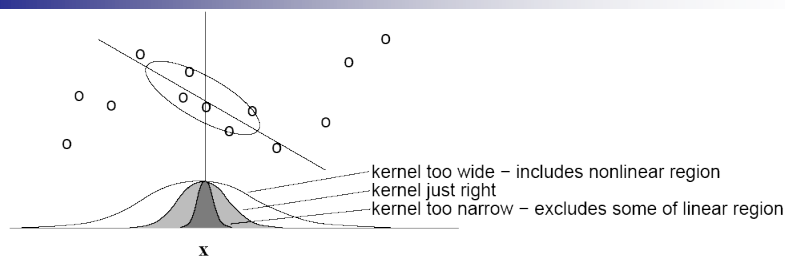
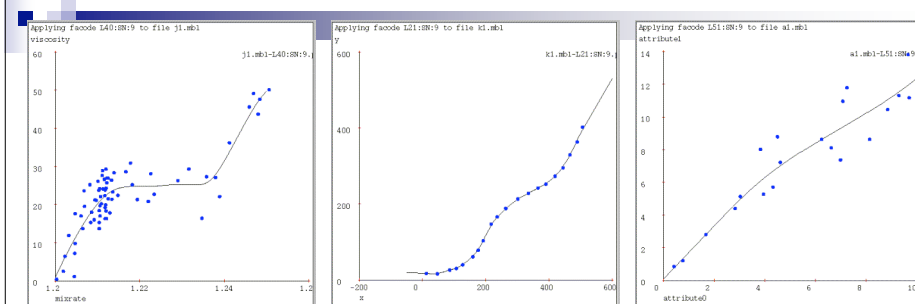


Image from Cohn, D.A., Ghahramani, Z., and Jordan, M.I. (1996) "Active Learning with Statistical Models", JAIR Volume 4, pages 149-145.

LWR on our test cases



KW = 1/16 of x-axis width.

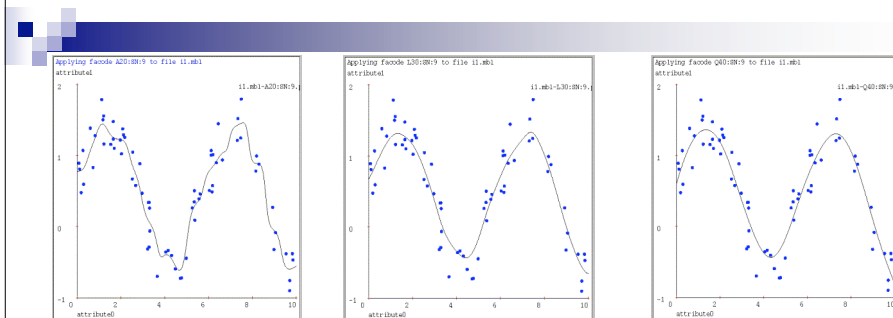
KW = 1/32 of x-axis width.

KW = 1/8 of x-axis width.

©2005-2007 Carlos Guestrin

45

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

46

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

47

Curse of the irrelevant feature

©2005-2007 Carlos Guestrin

48

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

49

Acknowledgment

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

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

©2005-2007 Carlos Guestrin

50