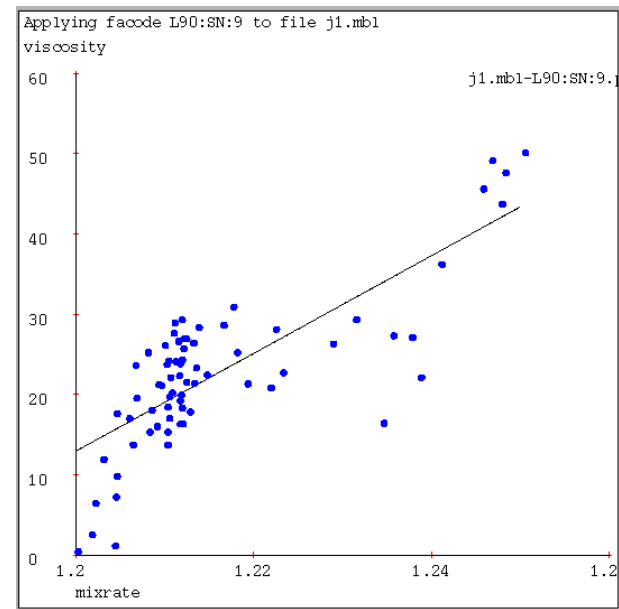
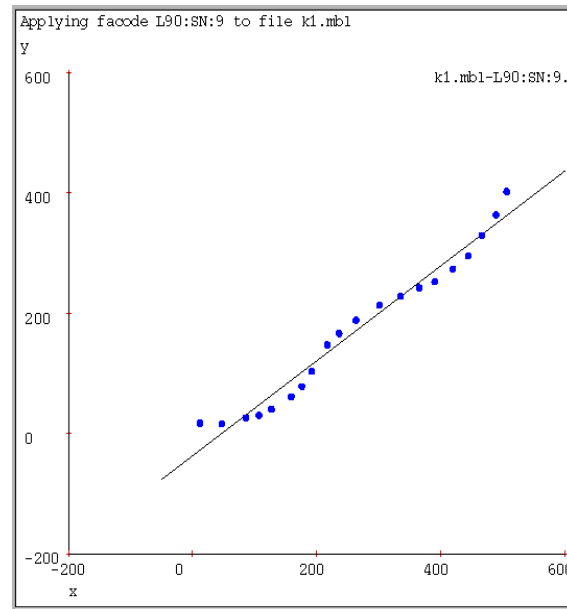
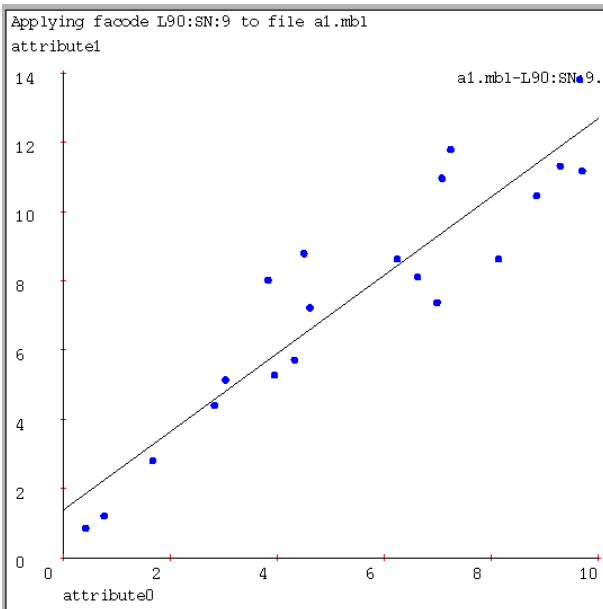


# CSE546: Instance-based Learning (a.k.a. non-parametric methods) Winter 2012

Luke Zettlemoyer

Slides adapted from Carlos Guestrin

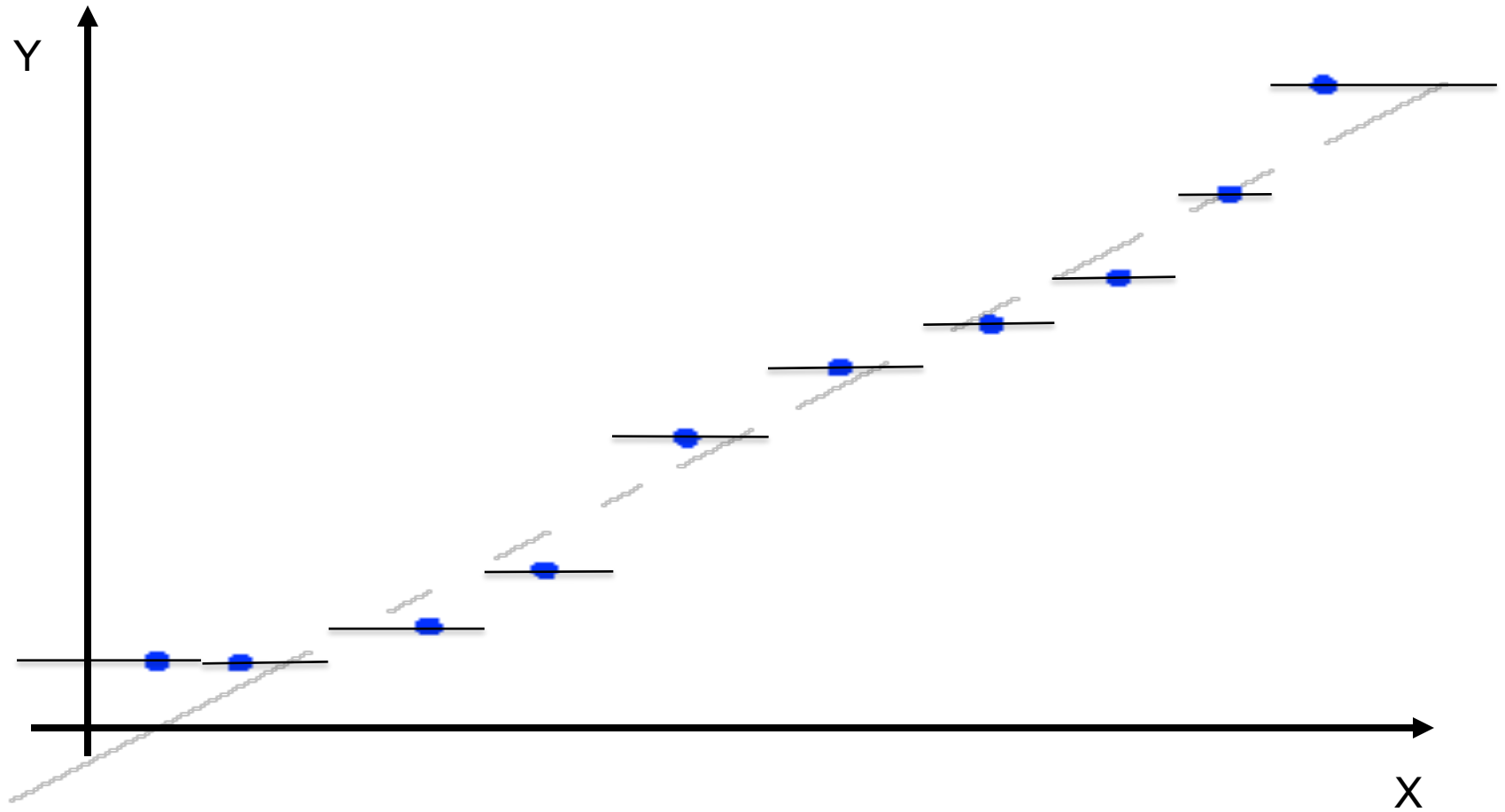
# Linear Regression: What can go wrong?



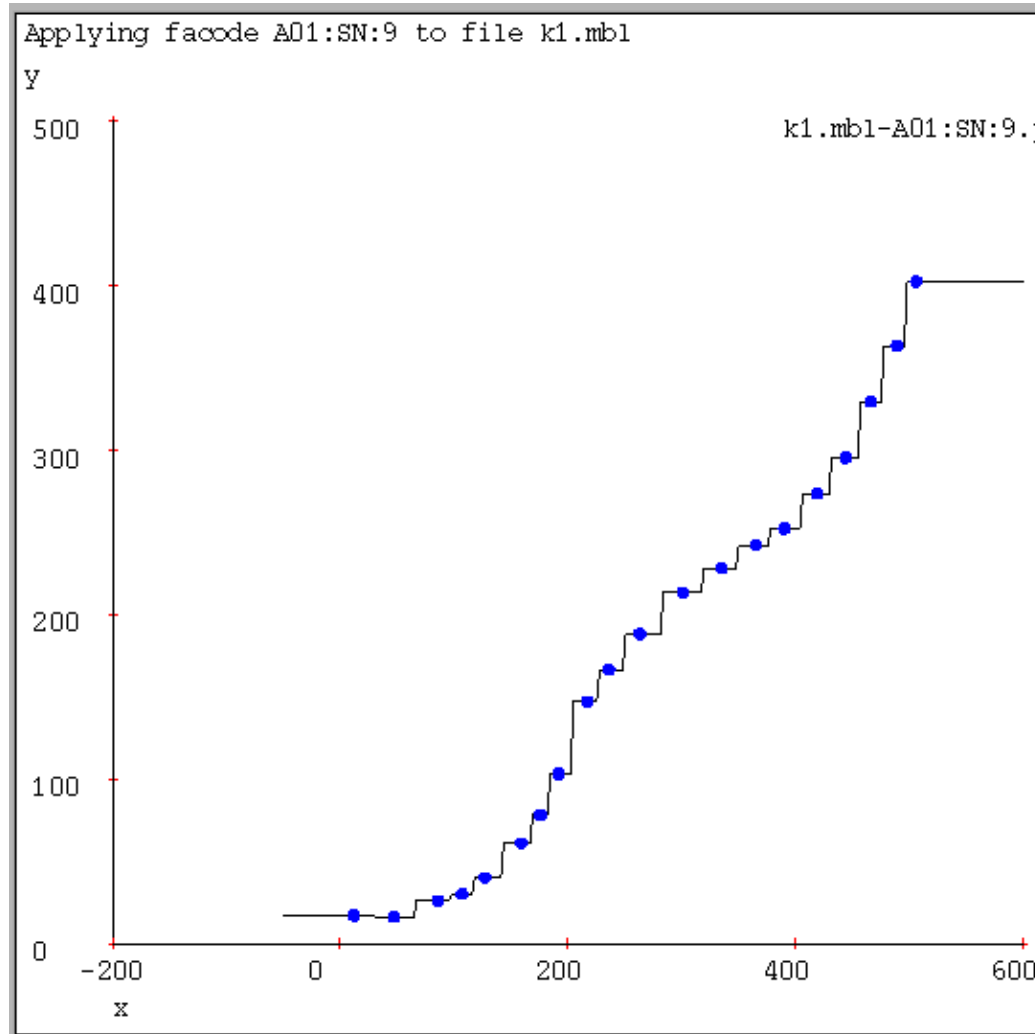
## What do we do if the bias is too strong?

- Might want the data to drive the complexity of the model!
- Try instance-based Learning (a.k.a. non-parametric methods)?

# Using data to predict new data



# Nearest neighbor with lots of data!



# Univariate 1-Nearest Neighbor

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

Given query point  $x_q$ , your job is to predict

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

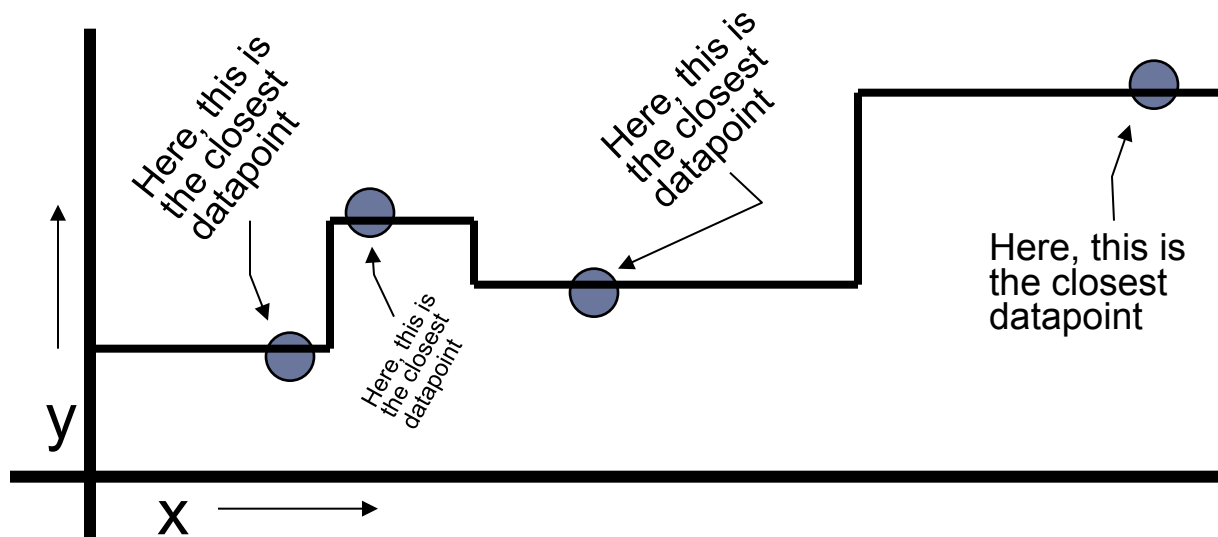
Nearest Neighbor:

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.

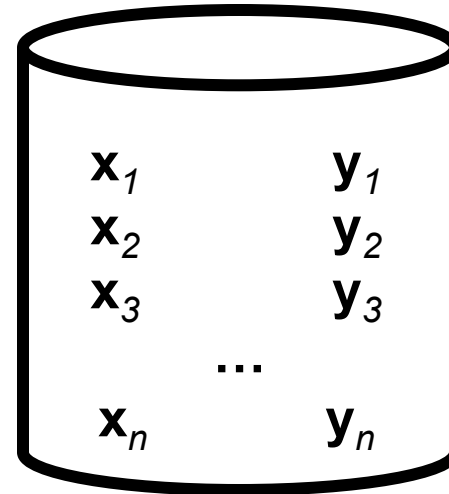


*1-Nearest Neighbor is an example of....*

## **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.



### **Instance-based learning, four things to specify:**

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

# 1-Nearest Neighbor

**Instance-based learning, four things to specify:**

1. *A distance metric*

**Often Euclidian (many more are possible)**

2. *How many nearby neighbors to look at?*

**One**

3. *A weighting function (optional)*

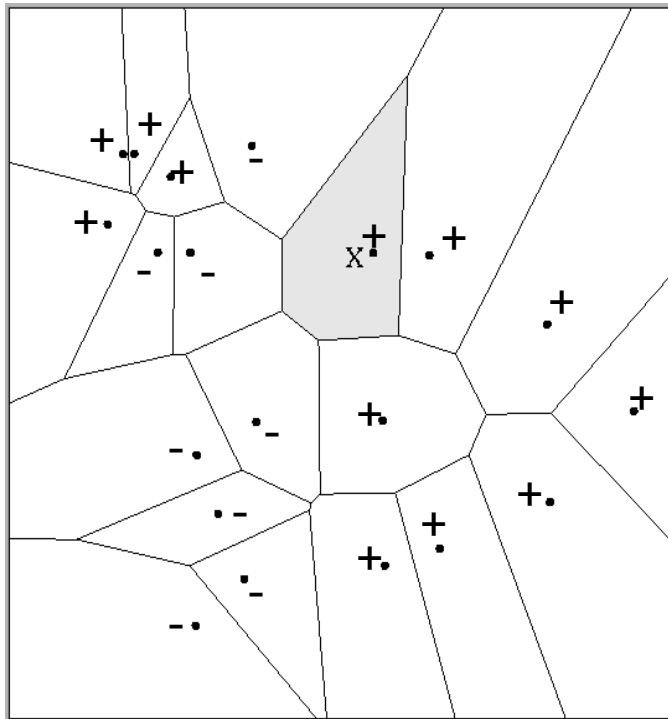
**Unused**

4. *How to fit with the local points?*

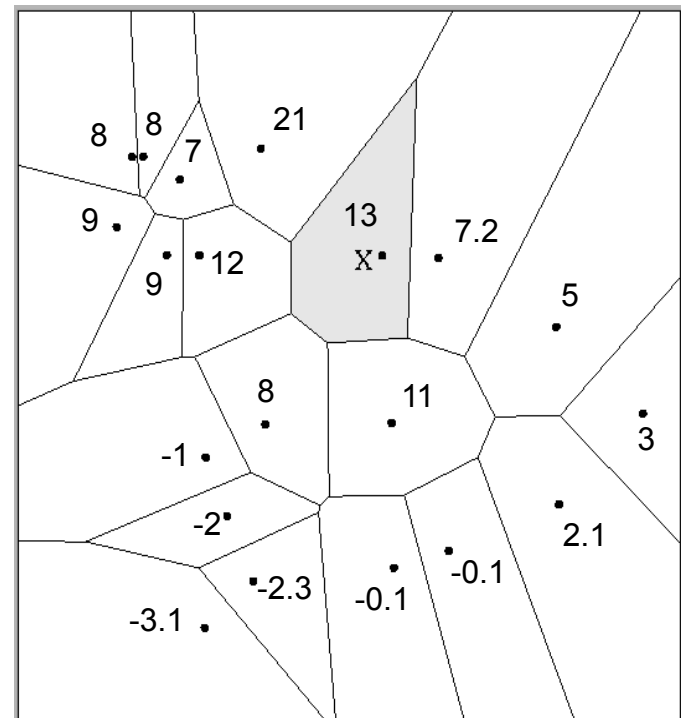
**Just predict the same output as the nearest neighbor.**

# Multivariate 1-NN examples

Classification

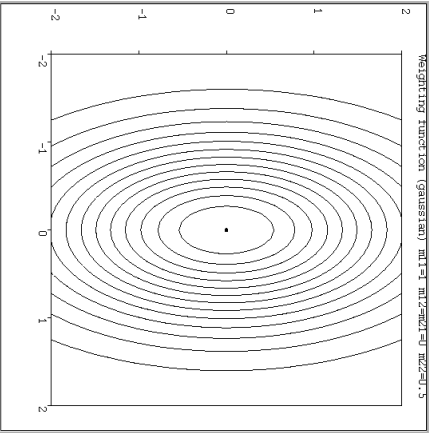


Regression

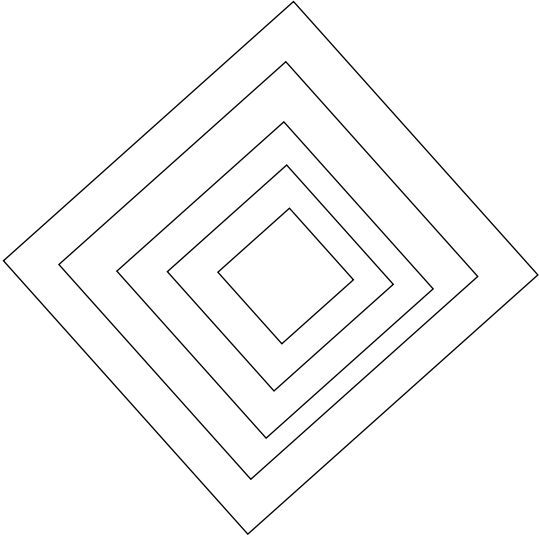




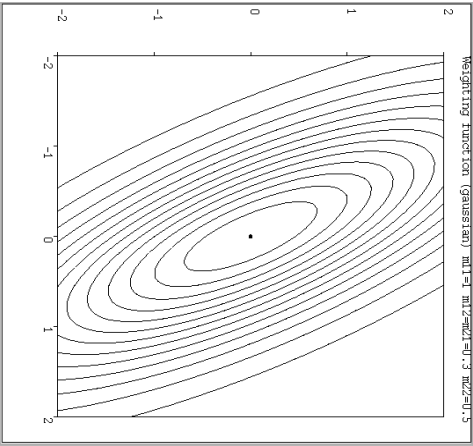
# Notable distance metrics (and their level sets)



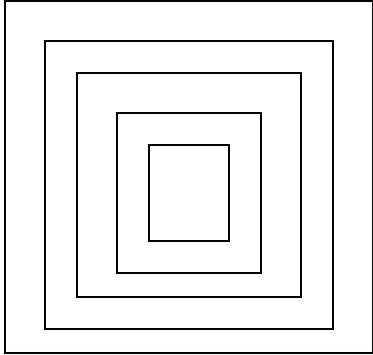
**Weighted Euclidian ( $L_2$ )**



**$L_1$  norm (absolute)**



**Mahalanobis**



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

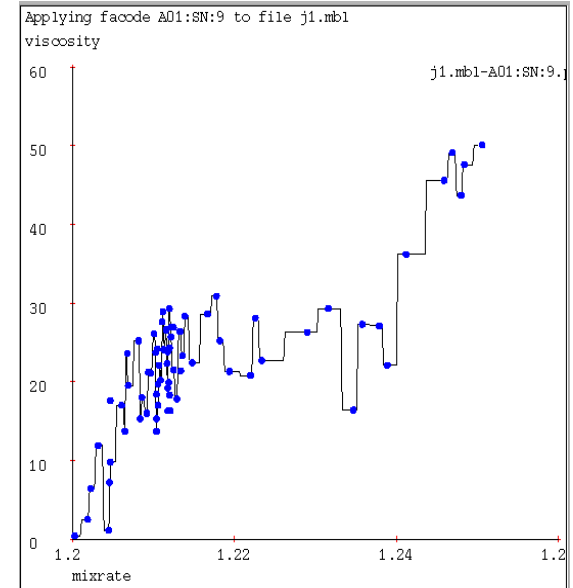
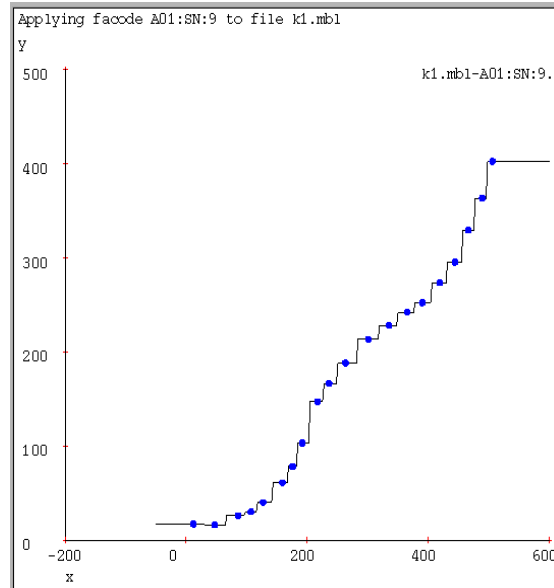
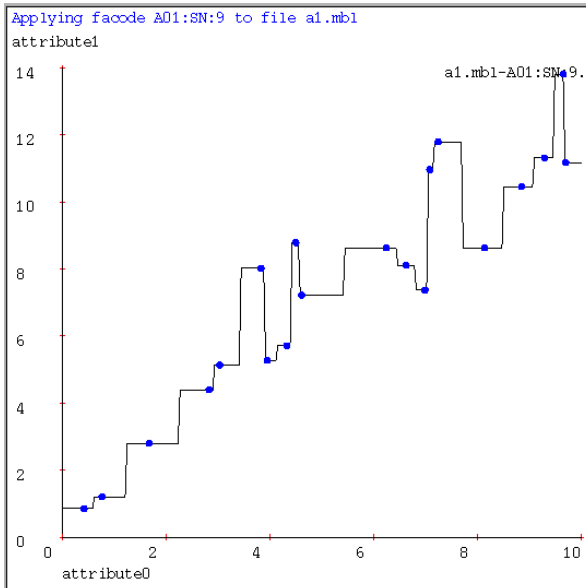
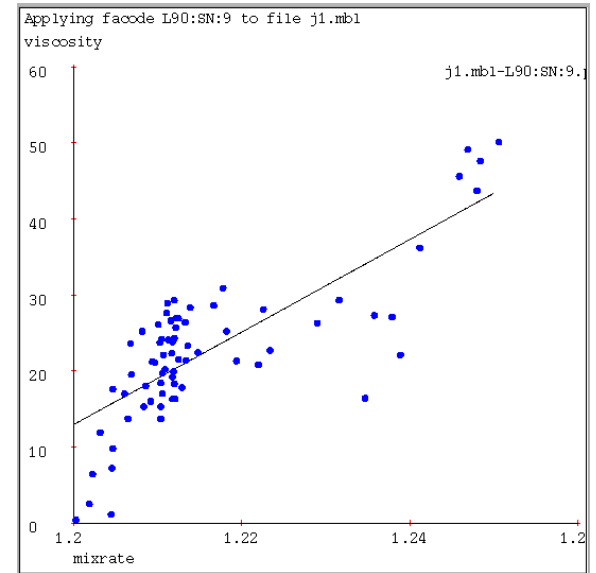
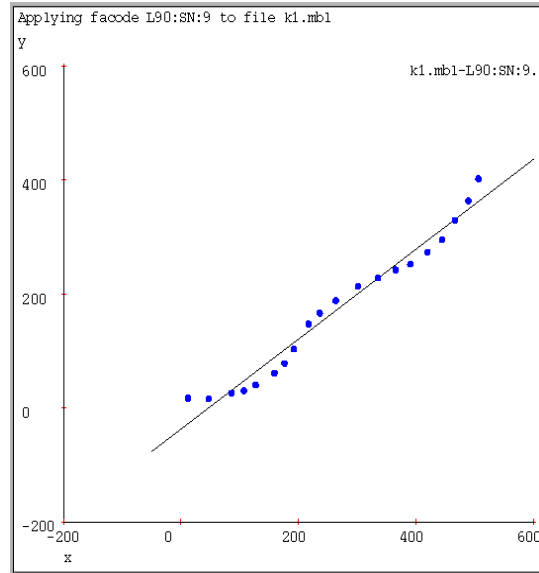
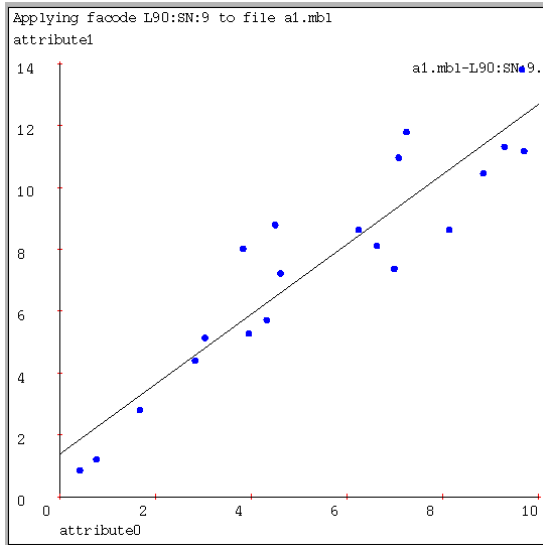
# Consistency of 1-NN

- 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 for any data distribution  $p(x)$ :

$$\lim_{n \rightarrow \infty} MSE(f_n) = 0 \quad MSE(f_n) = \int_x p(x) (f_n(x) - y_x)^2 dx$$

- **Linear regression is not consistent!**
  - Representation bias
- **1-NN is consistent**
  - What about noisy data?
  - What about variance?

# 1-NN overfits?



# k-Nearest Neighbor

**Instance-based learning, four things to specify:**

1. *A distance metric*

**Euclidian (and many more)**

2. *How many nearby neighbors to look at?*

**k**

1. *A weighting function (optional)*

**Unused**

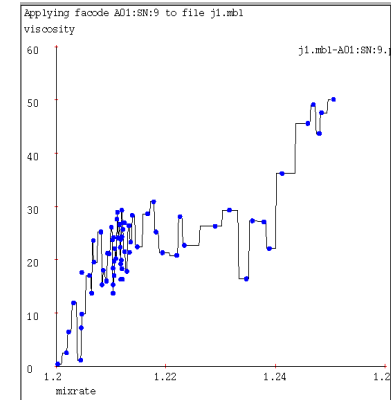
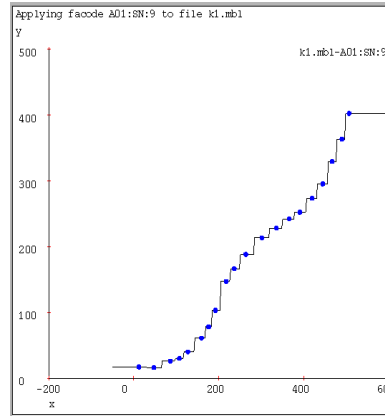
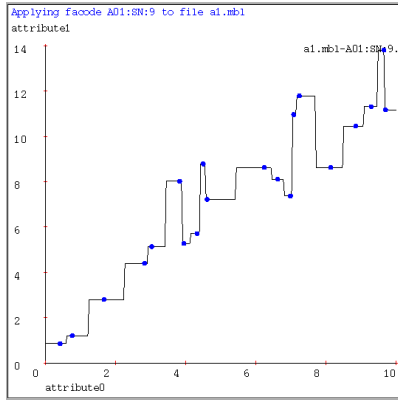
2. *How to fit with the local points?*

**Return the average output**

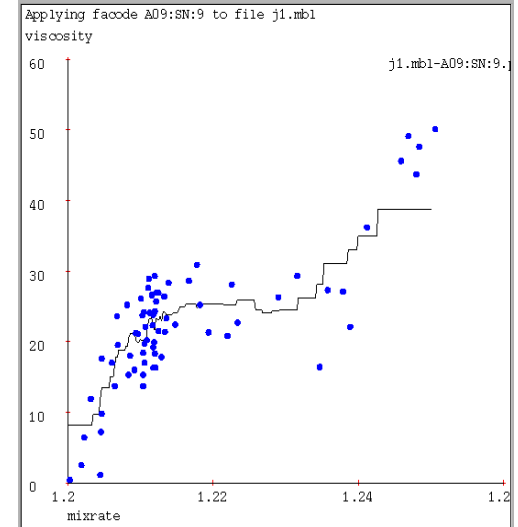
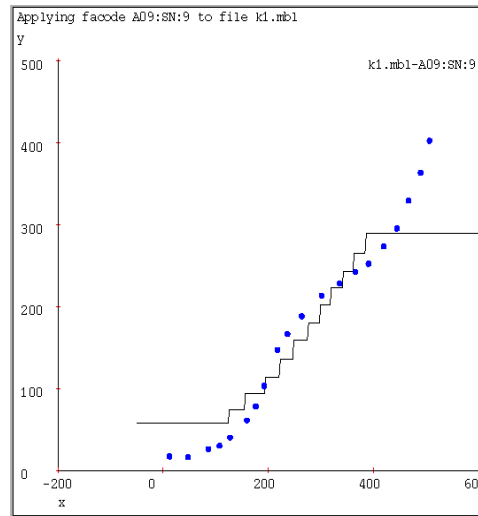
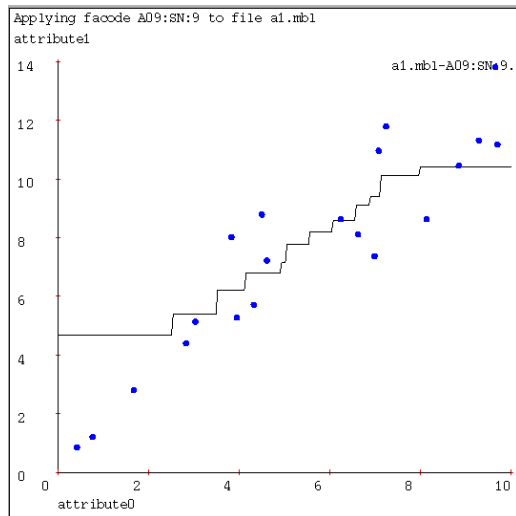
**predict:  $(1/k) \sum y_i$**  (summing over k nearest neighbors)

# k-Nearest Neighbor

k=1



k=9



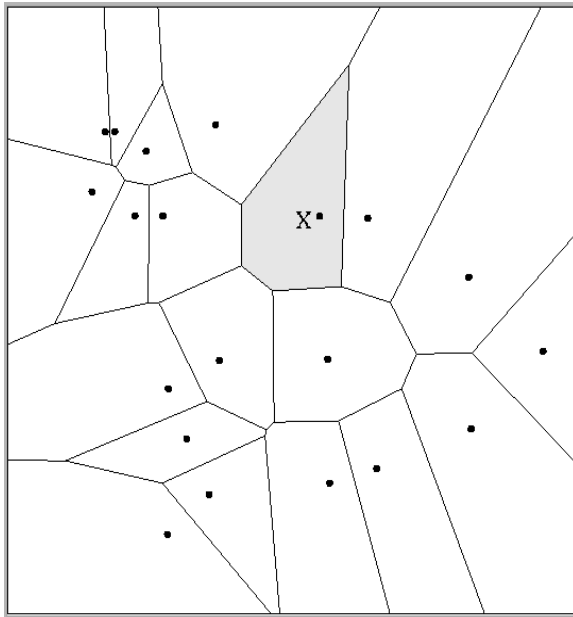
**Which is better? What can we do about the discontinuities?**

# Weighted distance metrics

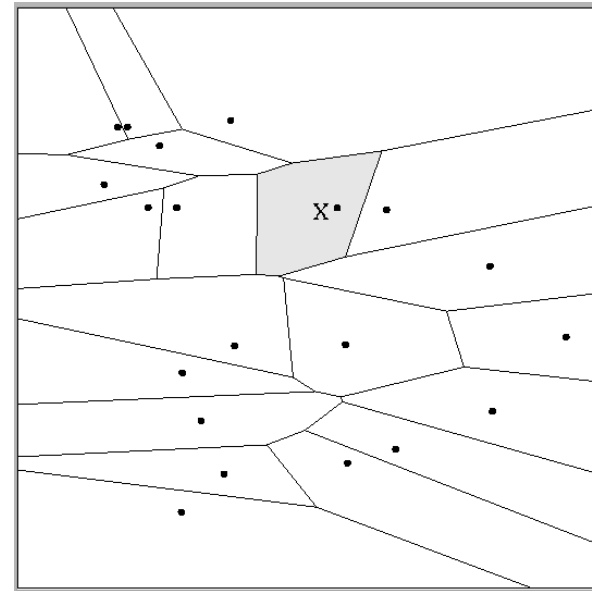
Suppose the input vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  are two dimensional:

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

Nearest-neighbor regions in input space:



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



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

The relative scaling of the distance metric affect region shapes

# Weighted Euclidean distance metric

Or equivalently,

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

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

where

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

Other Metrics...

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

# Kernel regression

## Instance-based learning:

1. *A distance metric*  
Euclidian (and many more)
2. *How many nearby neighbors to look at?*

All of them

3. *A weighting function*

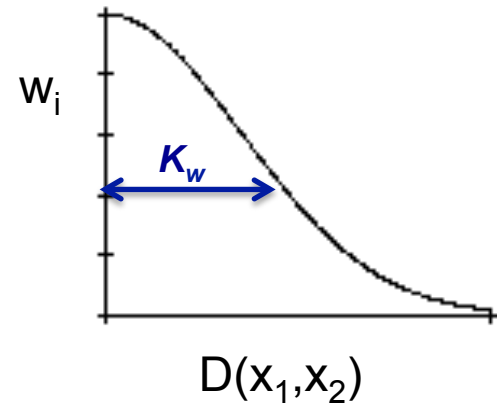
$$w_i = \exp(-D(x_i, \text{query})^2 / K_w^2)$$

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:

$$\text{predict} = \frac{\sum w_i y_i}{\sum w_i}$$



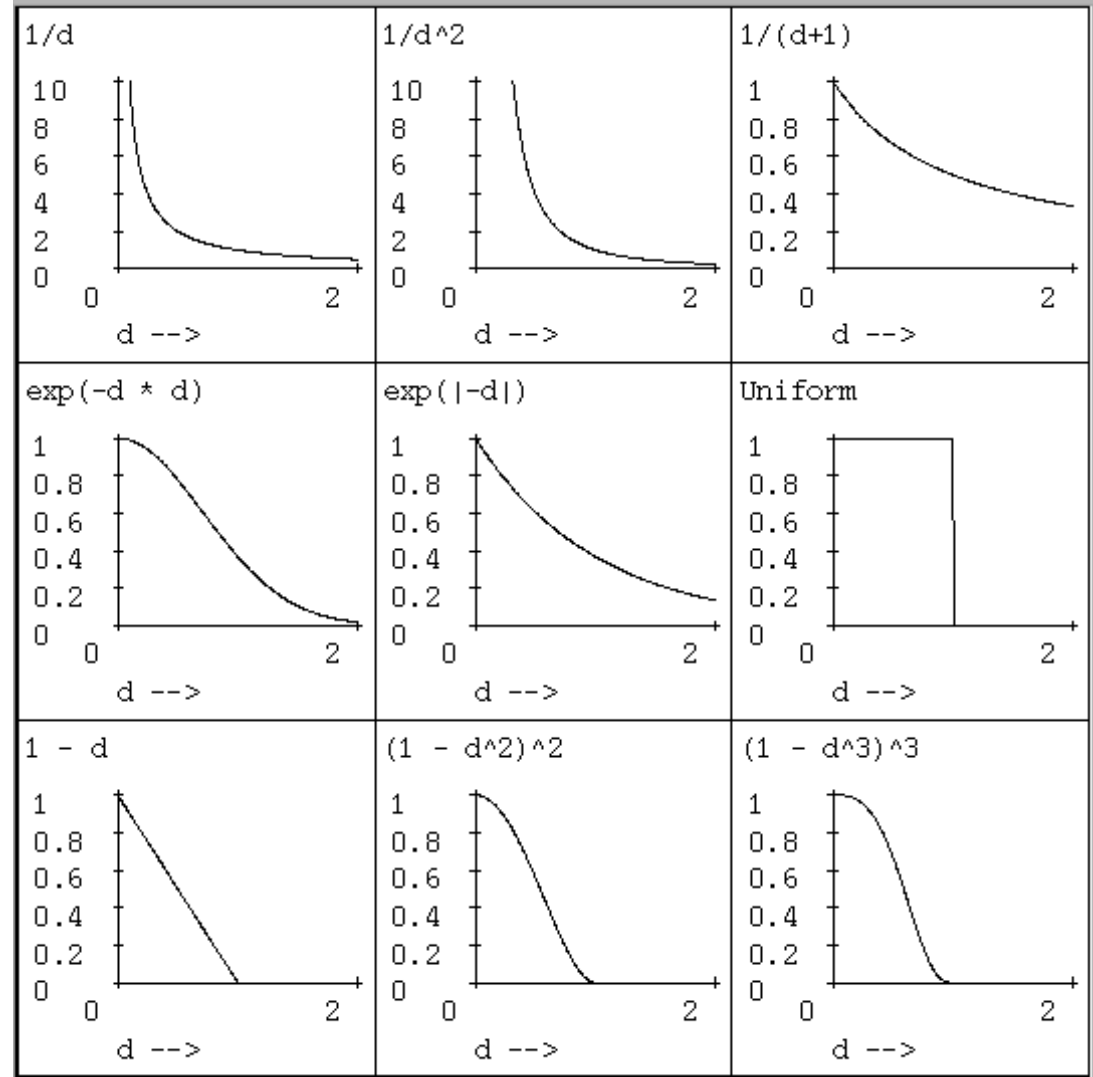


# Many possible weighting functions

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

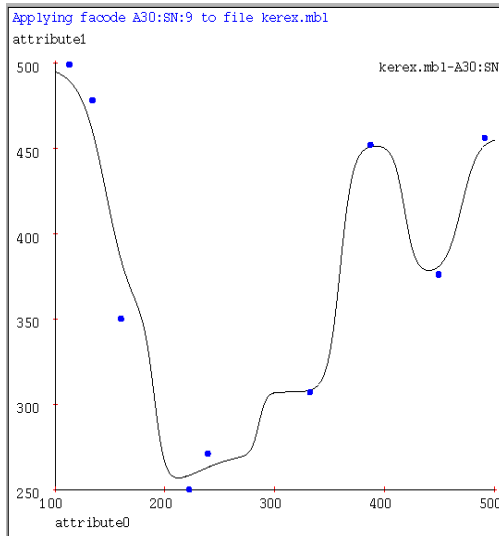
Typically:

- Choose D manually
- Optimize  $K_w$  using gradient descent

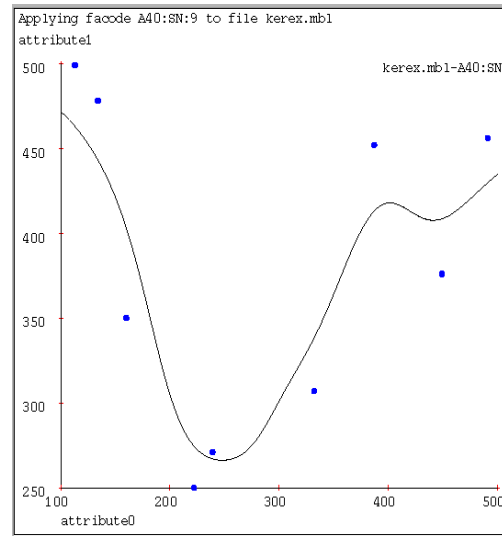


(Our examples use Gaussian)

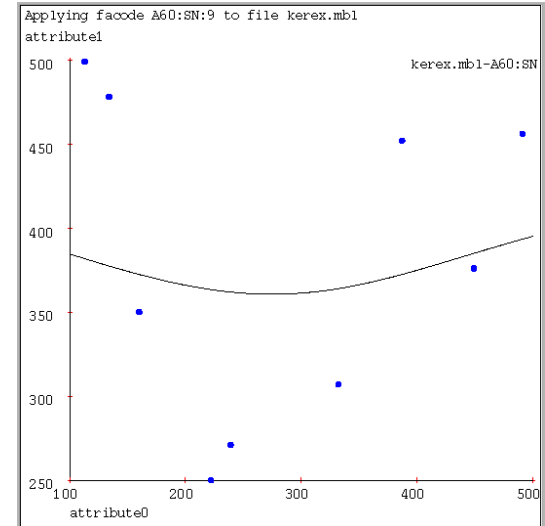
# Kernel regression predictions



$K_W=10$



$K_W=20$



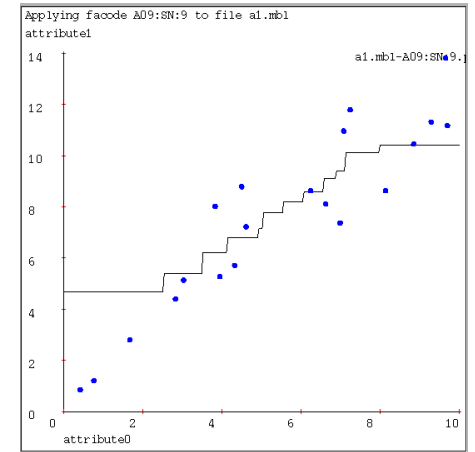
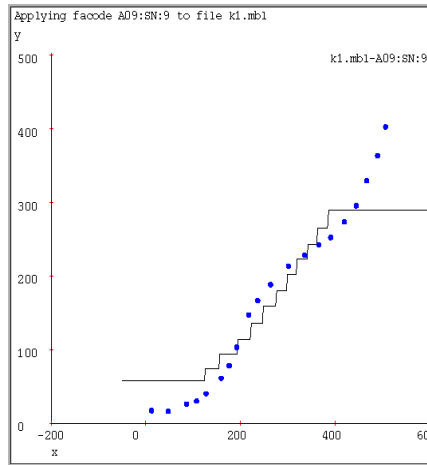
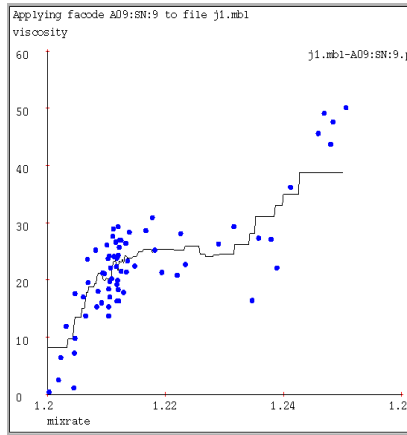
$K_W=80$

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

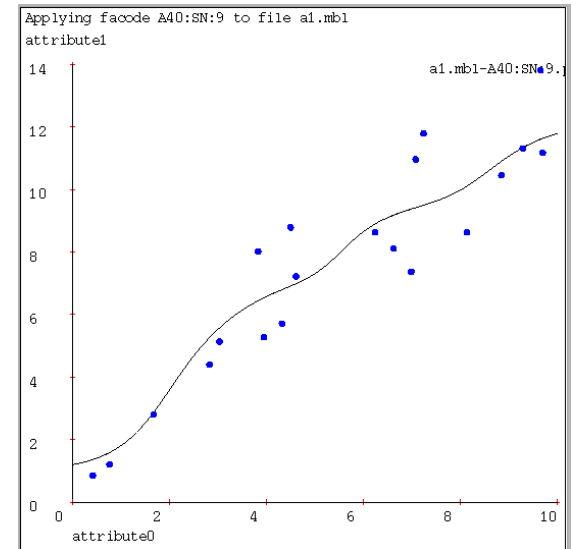
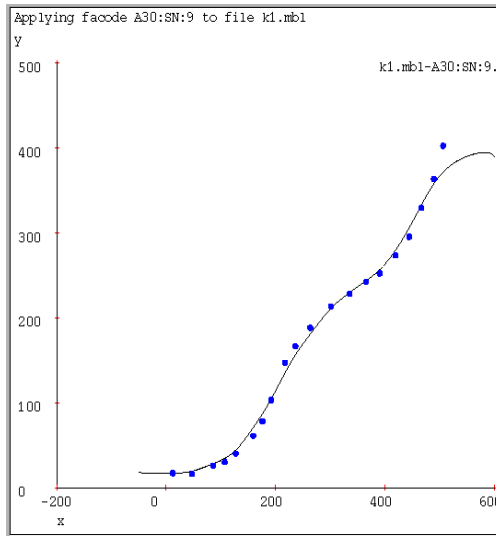
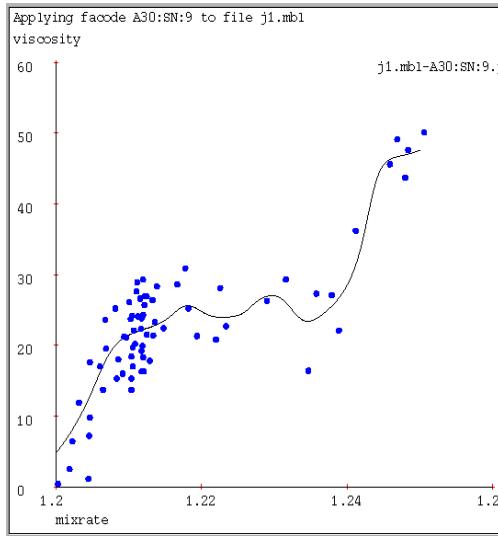
As  $K_W \rightarrow \infty$ , the prediction tends to the global average.

# Kernel regression on our test cases

NN k=9



Kernel regression



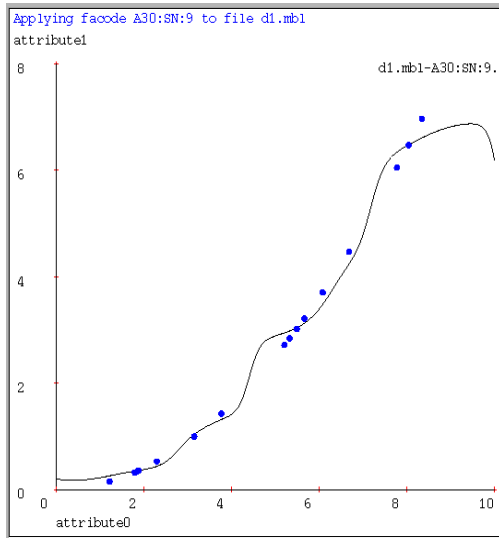
$K_W = 1/32$  of x-axis width.

$K_W = 1/32$  of x-axis width.

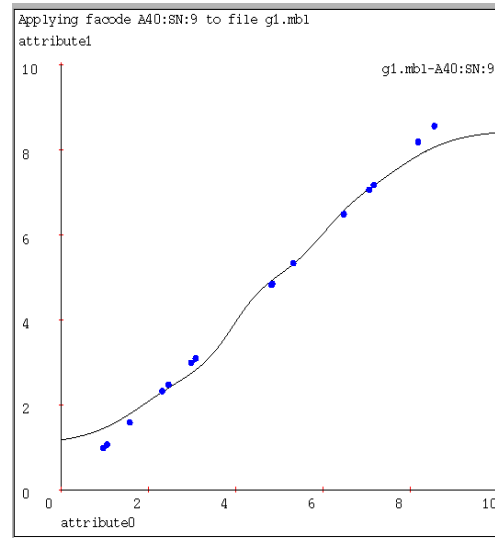
$K_W = 1/16$  axis width.

Choosing a good  $K_W$  is important! Remind you of anything we have seen?

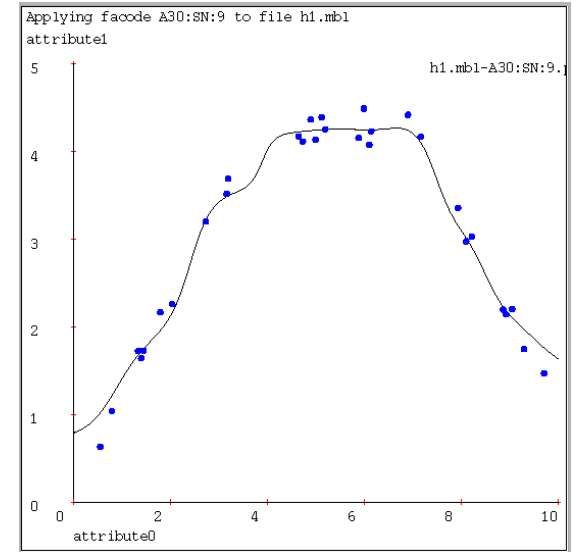
# Kernel regression: problem solved?



$K_W = \text{Best.}$



$K_W = \text{Best.}$



$K_W = \text{Best.}$

Where are we having problems?

- Sometimes in the middle...
- Generally, on the ends (extrapolation is hard!)

**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

## Instance-based learning, four things to specify:

- *A distance metric*

**Any**

- *How many nearby neighbors to look at?*

**All of them**

- *A weighting function (optional)*

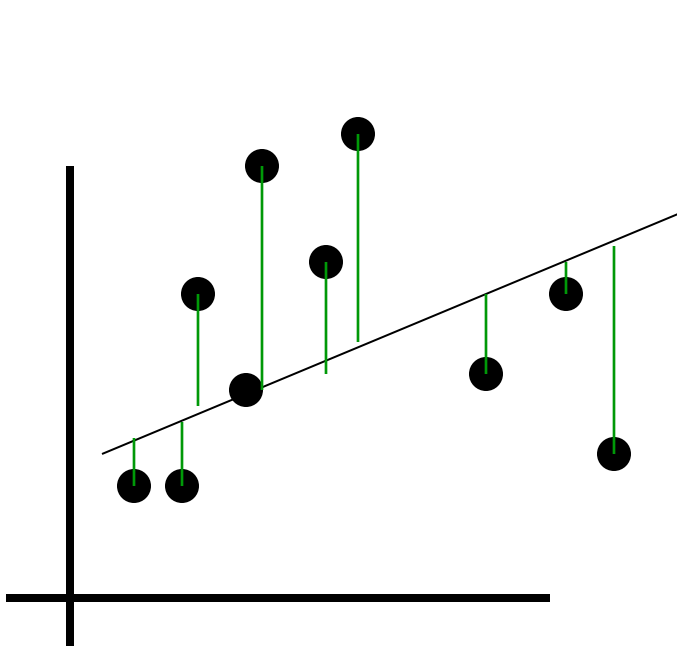
**Kernels:  $w_i = \exp(-D(x_i, query)^2 / Kw^2)$**

- *How to fit with the local points?*

**General weighted regression:**

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{k=1}^N w_k^2 (y_k - \beta^T \mathbf{x}_k)^2$$

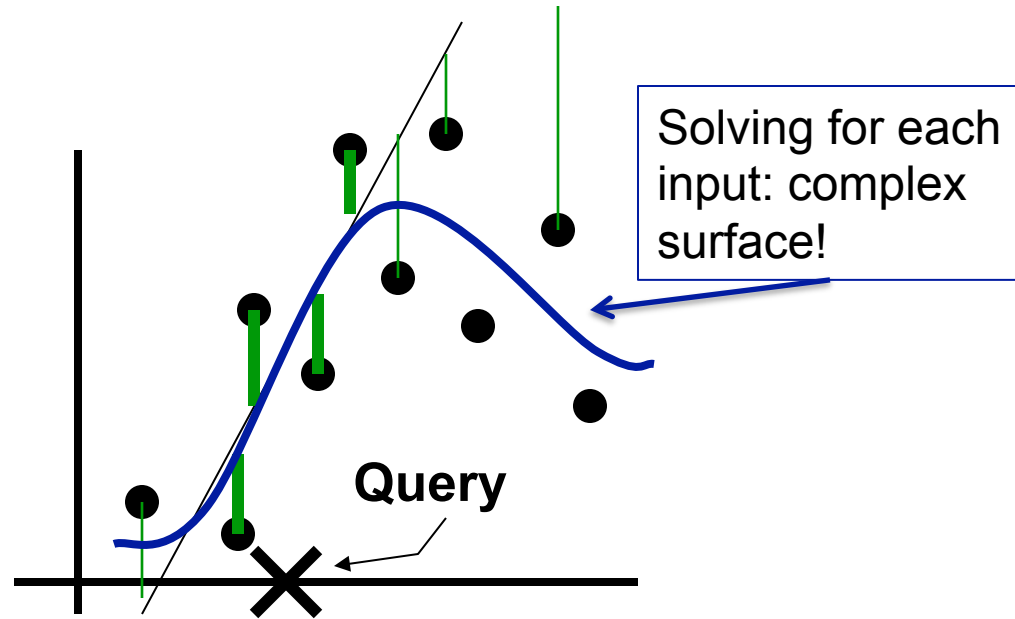
# How LWR works



## Linear regression

- Same parameters for all queries

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



Solving for each input: complex surface!

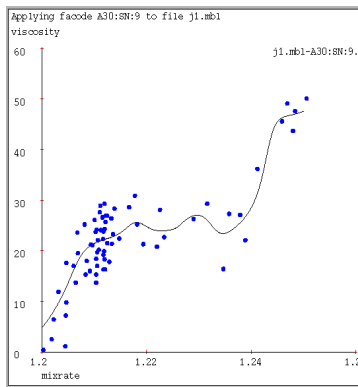
## Locally weighted regression

- Solve weighted linear regression for each query

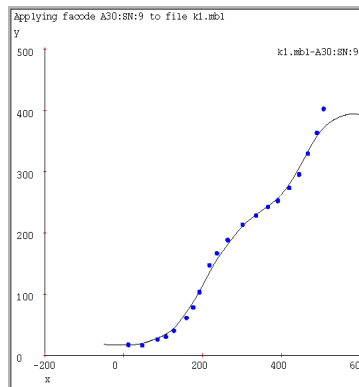
$$\beta = \left( (WX)^T WX \right)^{-1} (WX)^T WY$$
$$W = \begin{pmatrix} w_1 & 0 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & w_n \end{pmatrix}$$

# LWR on our test cases

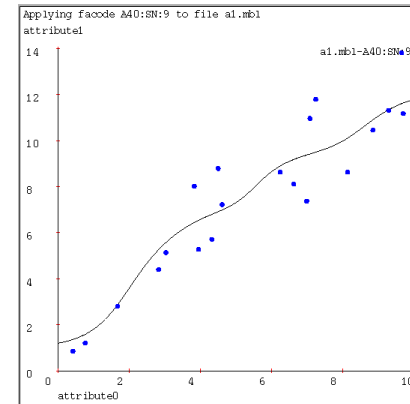
Kernel regression



$K_W = 1/32$  of x-axis width.

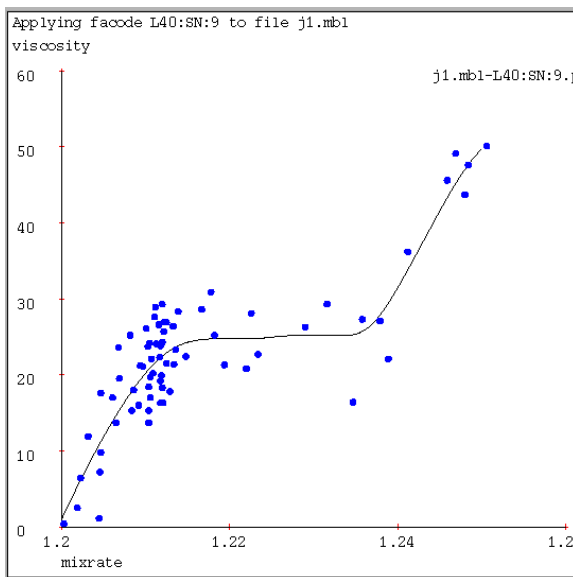


$K_W = 1/32$  of x-axis width.

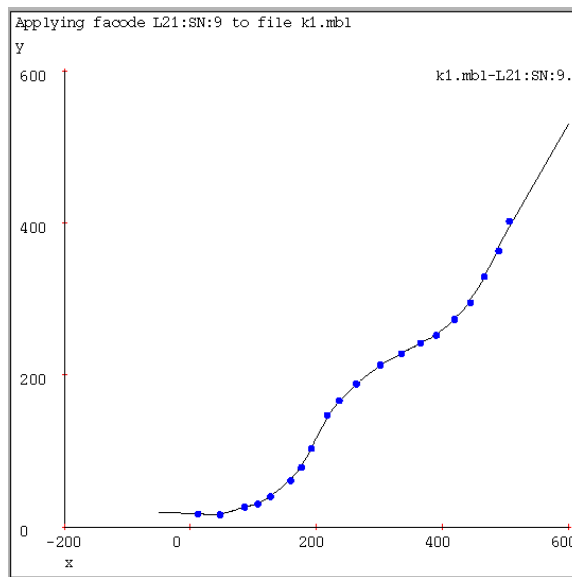


$K_W = 1/16$  axis width.

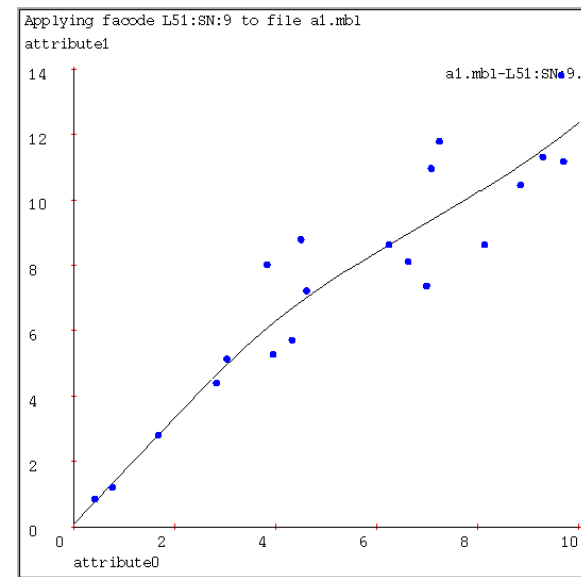
LWR



$K_W = 1/16$  of x-axis width.



$K_W = 1/32$  of x-axis width.

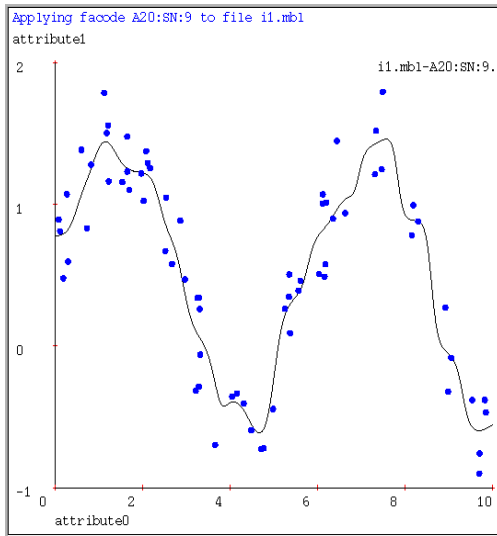


$K_W = 1/8$  of x-axis width.

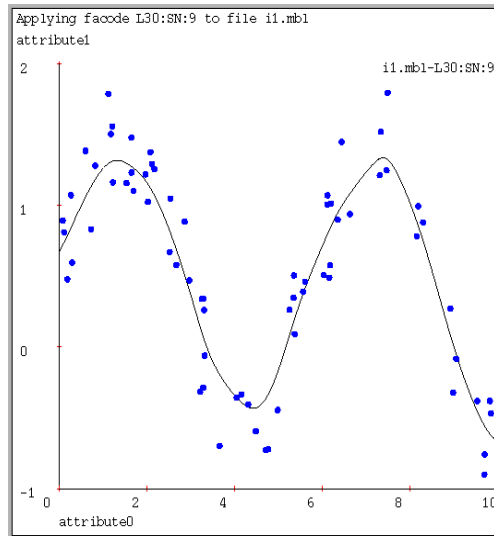


# Locally weighted polynomial regression

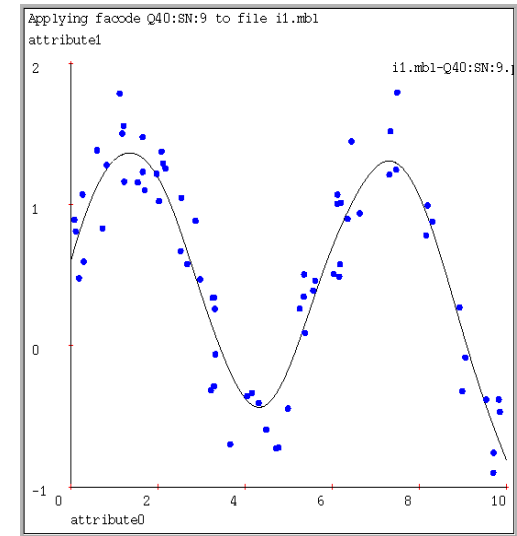
Kernel Regression: Kernel width  $K_W$  at optimal level.



$K_W = 1/100$  x-axis



$K_W = 1/40$  x-axis



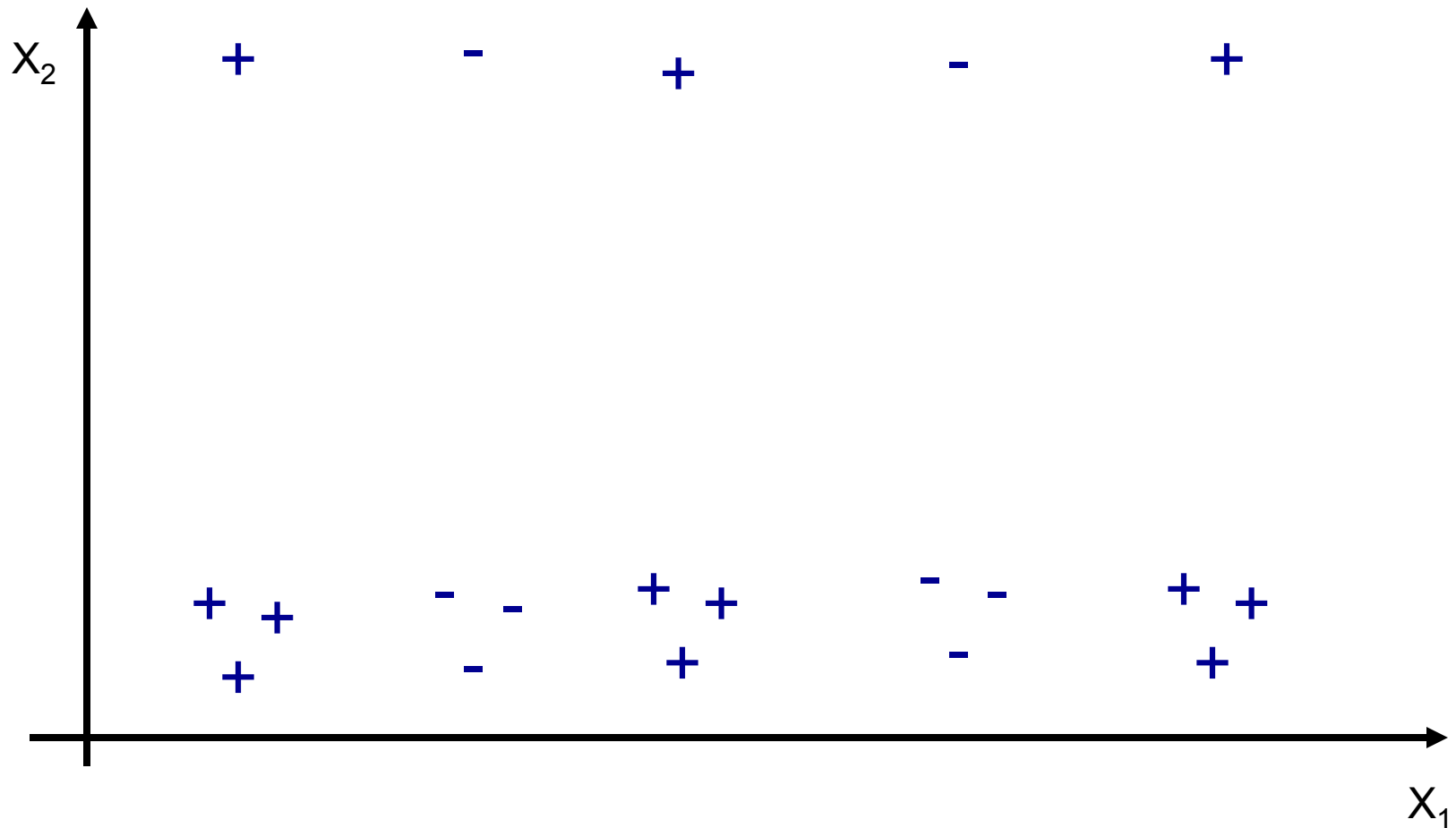
$K_W = 1/15$  x-axis

**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.

# Challenges for 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!
  - But, there are fast methods for dealing with large datasets
- **Instance-based learning often poor with noisy or irrelevant features**
  - In high dimensional spaces, all points will be very far from each other
  - Typically need a number of examples that is exponential in the dimension of  $X$
  - But, sometimes you are ok if you are clever about features

# Curse of the irrelevant feature



This is a contrived example, but similar problems are common in practice  
Need some form of feature selection!!

# 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>