

Predicting Real-valued outputs: an introduction to Regression

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials>. Comments and corrections gratefully received.

Andrew W. Moore
Professor
School of Computer Science
Carnegie Mellon University

www.cs.cmu.edu/~awm
awm@cs.cmu.edu
412-268-7599

This is reordered material from the Neural Nets lecture and the "Favorite Regression Algorithms" lecture

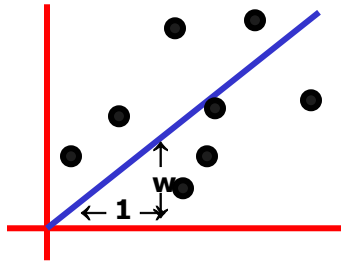
Copyright © 2001, 2003, Andrew W. Moore

Single-Parameter Linear Regression

Copyright © 2001, 2003, Andrew W. Moore

2

Linear Regression



DATASET

inputs	outputs
$x_1 = 1$	$y_1 = 1$
$x_2 = 3$	$y_2 = 2.2$
$x_3 = 2$	$y_3 = 2$
$x_4 = 1.5$	$y_4 = 1.9$
$x_5 = 4$	$y_5 = 3.1$

Linear regression assumes that the expected value of the output given an input, $E[y/x]$, is linear.

Simplest case: $\text{Out}(x) = wx$ for some unknown w .

Given the data, we can estimate w .

1-parameter linear regression

Assume that the data is formed by

$$y_i = wx_i + \text{noise}_i$$

where...

- the noise signals are independent
- the noise has a normal distribution with mean 0 and unknown variance σ^2

$p(y|w,x)$ has a normal distribution with

- mean wx
- variance σ^2

Bayesian Linear Regression

$$p(y|w, x) = \text{Normal}(\text{mean } wx, \text{var } \sigma^2)$$

We have a set of datapoints (x_1, y_1) (x_2, y_2) ... (x_n, y_n) which are **EVIDENCE** about w .

We want to infer w from the data.

$$p(w|x_1, x_2, x_3, \dots, x_n, y_1, y_2, \dots, y_n)$$

- You can use **BAYES** rule to work out a posterior distribution for w given the data.
- Or you could do Maximum Likelihood Estimation

Maximum likelihood estimation of w

Asks the question:

"For which value of w is this data most likely to have happened?"

\Leftrightarrow

For what w is

$$p(y_1, y_2, \dots, y_n | x_1, x_2, x_3, \dots, x_n, w) \text{ maximized?}$$

\Leftrightarrow

For what w is

$$\prod_{i=1}^n p(y_i | w, x_i) \text{ maximized'}$$

For what w is

$$\prod_{i=1}^n p(y_i | w, x_i) \text{ maximized?}$$

For what w is

$$\prod_{i=1}^n \exp\left(-\frac{1}{2} \left(\frac{y_i - wx_i}{\sigma}\right)^2\right) \text{ maximized?}$$

For what w is

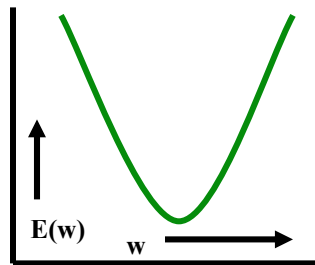
$$\sum_{i=1}^n -\frac{1}{2} \left(\frac{y_i - wx_i}{\sigma}\right)^2 \text{ maximized?}$$

For what w is

$$\sum_{i=1}^n (y_i - wx_i)^2 \text{ minimized?}$$

Linear Regression

The maximum likelihood w is the one that minimizes sum-of-squares of residuals



$$\begin{aligned} E &= \sum_i (y_i - wx_i)^2 \\ &= \sum_i y_i^2 - (2 \sum_i x_i y_i)w + (\sum_i x_i^2)w^2 \end{aligned}$$

We want to minimize a quadratic function of w .

Linear Regression

Easy to show the sum of squares is minimized when

$$w = \frac{\sum x_i y_i}{\sum x_i^2}$$

The maximum likelihood model is $\text{Out}(x) = wx$

We can use it for prediction

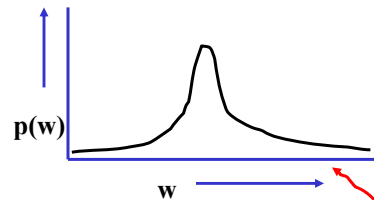
Linear Regression

Easy to show the sum of squares is minimized when

$$w = \frac{\sum x_i y_i}{\sum x_i^2}$$

The maximum likelihood model is $\text{Out}(x) = wx$

We can use it for prediction



Note: In Bayesian stats you'd have ended up with a prob dist of w

And predictions would have given a prob dist of expected output

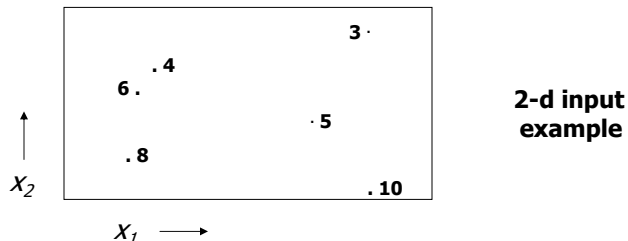
Often useful to know your confidence.

Max likelihood can give some kinds of confidence too.

Multivariate Linear Regression

Multivariate Regression

What if the inputs are vectors?



Dataset has form

$$\begin{array}{ll} \mathbf{x}_1 & y_1 \\ \mathbf{x}_2 & y_2 \\ \mathbf{x}_3 & y_3 \\ \vdots & \vdots \\ \mathbf{x}_R & y_R \end{array}$$

Multivariate Regression

Write matrix X and Y thus:

$$\mathbf{X} = \begin{bmatrix} \dots \mathbf{X}_1 \dots \\ \dots \mathbf{X}_2 \dots \\ \vdots \\ \dots \mathbf{X}_R \dots \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{R1} & x_{R2} & \dots & x_{Rm} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_R \end{bmatrix}$$

(there are R datapoints. Each input has m components)

The linear regression model assumes a vector \mathbf{w} such that

$$\text{Out}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = w_1 x[1] + w_2 x[2] + \dots w_m x[D]$$

The max. likelihood \mathbf{w} is $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y})$

Multivariate Regression

Write matrix X and Y thus:

$$\mathbf{X} = \begin{bmatrix} \dots \mathbf{X}_1 \dots \\ \dots \mathbf{X}_2 \dots \\ \vdots \\ \dots \mathbf{X}_R \dots \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{R1} & x_{R2} & \dots & x_{Rm} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_R \end{bmatrix}$$

(there are R datapoints. Each input

**IMPORTANT EXERCISE:
PROVE IT !!!!!**

The linear regression model assumes a vector \mathbf{w} such that

$$\text{Out}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = w_1 x[1] + w_2 x[2] + \dots w_m x[D]$$

The max. likelihood \mathbf{w} is $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y})$

Multivariate Regression (con't)

The max. likelihood \mathbf{w} is $\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{Y})$

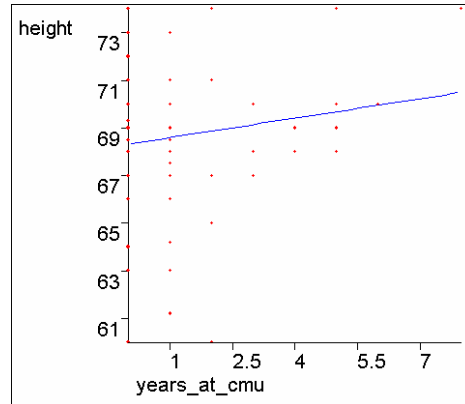
$\mathbf{X}^T\mathbf{X}$ is an $m \times m$ matrix: i, j 'th elt is $\sum_{k=1}^R x_{ki}x_{kj}$

$\mathbf{X}^T\mathbf{Y}$ is an m -element vector: i 'th elt $\sum_{k=1}^R x_{ki}y_k$

Constant Term in Linear Regression

What about a constant term?

We may expect linear data that does not go through the origin.



Statisticians and Neural Net Folks all agree on a simple obvious hack.

Can you guess??

The constant term

- The trick is to create a fake input " X_0 " that always takes the value 1

X_1	X_2	Y
2	4	16
3	4	17
5	5	20

Before:

$$Y = w_1 X_1 + w_2 X_2$$

...has to be a poor model

In this example, You should be able to see the MLE w_0 , w_1 and w_2 by inspection

X_0	X_1	X_2	Y
1	2	4	16
1	3	4	17
1	5	5	20

After:

$$Y = w_0 X_0 + w_1 X_1 + w_2 X_2$$

$$= w_0 + w_1 X_1 + w_2 X_2$$

...has a fine constant term

Heteroscedasticity...

Linear Regression with varying noise

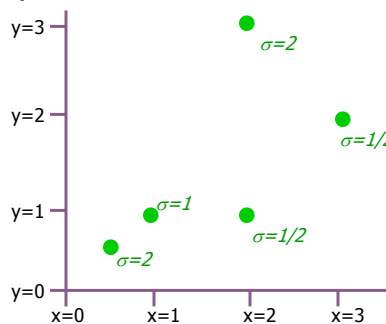
Copyright © 2001, 2003, Andrew W. Moore

19

Regression with varying noise

- Suppose you know the variance of the noise that was added to each datapoint.

x_i	y_i	σ_i^2
$1/2$	$1/2$	4
1	1	1
2	1	$1/4$
2	3	4
3	2	$1/4$



Assume $y_i \sim N(wx_i, \sigma_i^2)$

What's the MLE estimate of w ?

Copyright © 2001, 2003, Andrew W. Moore

20

MLE estimation with varying noise

$$\operatorname{argmax}_w \log p(y_1, y_2, \dots, y_R | x_1, x_2, \dots, x_R, \sigma_1^2, \sigma_2^2, \dots, \sigma_R^2, w) =$$

w

$$\operatorname{argmin}_w \sum_{i=1}^R \frac{(y_i - wx_i)^2}{\sigma_i^2} =$$

Assuming independence among noise and then plugging in equation for Gaussian and simplifying.

$$\left(w \text{ such that } \sum_{i=1}^R \frac{x_i(y_i - wx_i)}{\sigma_i^2} = 0 \right) =$$

Setting dLL/dw equal to zero

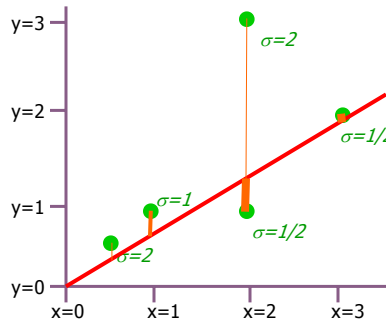
$$\frac{\left(\sum_{i=1}^R \frac{x_i y_i}{\sigma_i^2} \right)}{\left(\sum_{i=1}^R \frac{x_i^2}{\sigma_i^2} \right)}$$

Trivial algebra

This is Weighted Regression

- We are asking to minimize the weighted sum of squares

$$\operatorname{argmin}_w \sum_{i=1}^R \frac{(y_i - wx_i)^2}{\sigma_i^2}$$



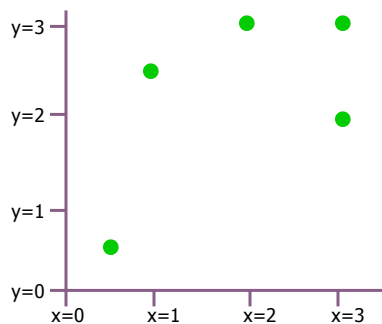
where weight for i 'th datapoint is $\frac{1}{\sigma_i^2}$

Non-linear Regression

Non-linear Regression

- Suppose you know that y is related to a function of x in such a way that the predicted values have a non-linear dependence on w , e.g:

x_i	y_i
$1/2$	$1/2$
1	2.5
2	3
3	2
3	3



Assume $y_i \sim N(\sqrt{w + x_i}, \sigma^2)$

What's the MLE estimate of w ?

Non-linear MLE estimation

$$\operatorname{argmax}_w \log p(y_1, y_2, \dots, y_R \mid x_1, x_2, \dots, x_R, \sigma, w) =$$

w

$$\operatorname{argmin}_w \sum_{i=1}^R (y_i - \sqrt{w + x_i})^2 =$$

Assuming i.i.d. and then plugging in equation for Gaussian and simplifying.

$$\left(w \text{ such that } \sum_{i=1}^R \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0 \right) =$$

Setting dLL/dw equal to zero

Non-linear MLE estimation

$$\operatorname{argmax}_w \log p(y_1, y_2, \dots, y_R \mid x_1, x_2, \dots, x_R, \sigma, w) =$$

w

$$\operatorname{argmin}_w \sum_{i=1}^R (y_i - \sqrt{w + x_i})^2 =$$

Assuming i.i.d. and then plugging in equation for Gaussian and simplifying.

$$\left(w \text{ such that } \sum_{i=1}^R \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0 \right) =$$

Setting dLL/dw equal to zero



We're down the algebraic toilet

So guess what we do?

Non-linear MLE estimation

$$\operatorname{argmax}_w \log p(y_1, y_2, \dots, y_R \mid x_1, x_2, \dots, x_R, \sigma, w) =$$

Common (but not only) approach:
Numerical Solutions:

- Line Search
- Simulated Annealing
- Gradient Descent
- Conjugate Gradient
- Levenberg Marquart
- Newton's Method

$$\left(y + x_i \right)^2 =$$

Assuming i.i.d. and then plugging in equation for Gaussian and simplifying.

$$\left(\frac{\partial}{\partial w} \left(y + x_i \right)^2 = 0 \right) =$$

Setting dLL/dw equal to zero

We're down the algebraic toilet

Also, special purpose statistical-optimization-specific tricks such as E.M. (See Gaussian Mixtures lecture for introduction)



So guess what we do?

Polynomial Regression

Polynomial Regression

So far we've mainly been dealing with linear regression

X_1	X_2	Y
3	2	7
1	1	3
⋮	⋮	⋮

x	y
3	7
2	3
1	1
⋮	⋮

$y_1=7..$

z	y
1	7
3	3
2	1
⋮	⋮

$z_1=(1,3,2).. \quad y_1=7..$
 $z_k=(1, x_{k1}, x_{k2})$

$$\beta = (Z^T Z)^{-1} (Z^T y)$$

$$y^{est} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

Quadratic Regression

It's trivial to do linear fits of fixed nonlinear basis functions

X_1	X_2	Y
3	2	7
1	1	3
⋮	⋮	⋮

x	y
3	7
2	3
1	1
⋮	⋮

$y_1=7..$

z	y
1	7
3	3
2	1
9	⋮
6	⋮
4	⋮
⋮	⋮

$z=(1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$

$$\beta = (Z^T Z)^{-1} (Z^T y)$$

$$y^{est} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_1 x_2 + \beta_5 x_2^2$$

Quadratic Regression

It's trivial

X_1	X_2
3	2
1	1
⋮	⋮

$$\mathbf{z} = \begin{pmatrix} 1 \\ 1 \\ \vdots \end{pmatrix}$$

$$\mathbf{z} = (1$$

Each component of a \mathbf{z} vector is called a term.
 Each column of the \mathbf{Z} matrix is called a term column
 How many terms in a quadratic regression with m inputs?
 • 1 constant term
 • m linear terms
 • $(m+1)\text{-choose-}2 = m(m+1)/2$ quadratic terms
 $(m+2)\text{-choose-}2$ terms in total = $O(m^2)$
 Note that solving $\beta = (\mathbf{Z}^T \mathbf{Z})^{-1} (\mathbf{Z}^T \mathbf{y})$ is thus $O(m^6)$

Qth-degree polynomial Regression

X_1	X_2	Y
3	2	7
1	1	3
⋮	⋮	⋮

$$\mathbf{x} = \begin{pmatrix} 3 & 2 \\ 1 & 1 \\ \vdots & \vdots \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 7 \\ 3 \\ \vdots \end{pmatrix}$$

$$\mathbf{z} = \begin{pmatrix} 1 & 3 & 2 & 9 & 6 & \dots \\ 1 & 1 & 1 & 1 & 1 & \dots \\ \vdots & & & & & \dots \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 7 \\ 3 \\ \vdots \end{pmatrix}$$

$\mathbf{z} =$ (all products of powers of inputs in which sum of powers is q or less)

$$\beta = (\mathbf{Z}^T \mathbf{Z})^{-1} (\mathbf{Z}^T \mathbf{y})$$

$$y^{est} = \beta_0 + \beta_1 x_1 + \dots$$

m inputs, degree Q: how many terms?

= the number of unique terms of the form

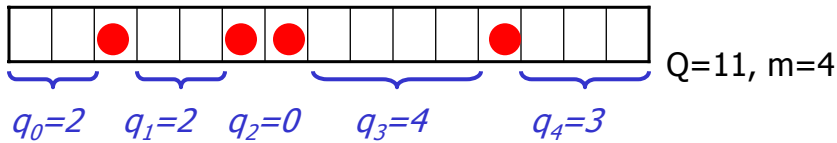
$$x_1^{q_1} x_2^{q_2} \dots x_m^{q_m} \text{ where } \sum_{i=1}^m q_i \leq Q$$

= the number of unique terms of the form

$$1^{q_0} x_1^{q_1} x_2^{q_2} \dots x_m^{q_m} \text{ where } \sum_{i=0}^m q_i = Q$$

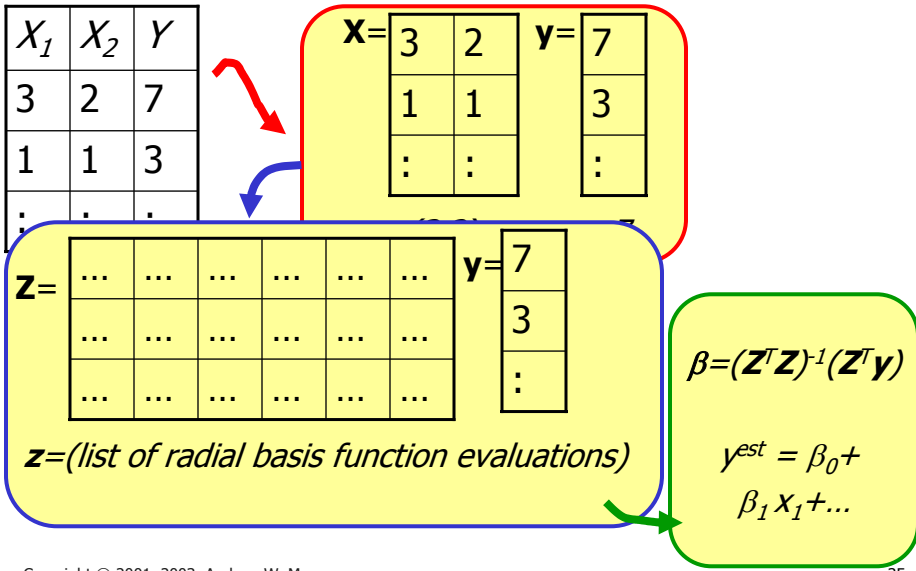
= the number of lists of non-negative integers $[q_0, q_1, q_2, \dots, q_m]$ in which $\sum q_i = Q$

= the number of ways of placing Q red disks on a row of squares of length Q+m = (Q+m)-choose-Q

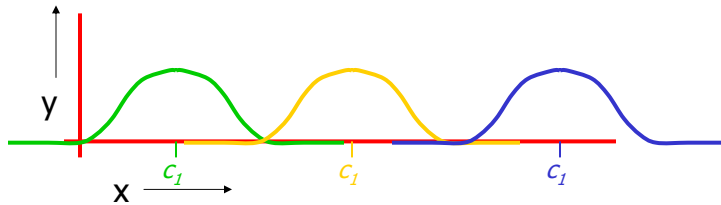


Radial Basis Functions

Radial Basis Functions (RBFs)



1-d RBFs

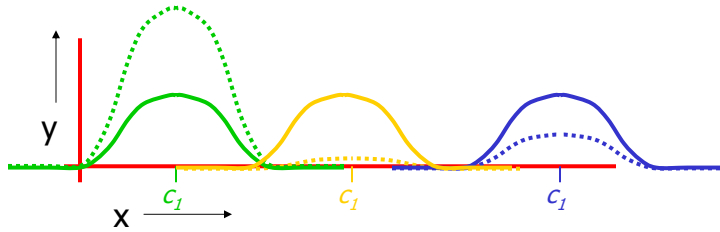


$$y^{est} = \beta_1 \phi_1(x) + \beta_2 \phi_2(x) + \beta_3 \phi_3(x)$$

where

$$\phi_i(x) = \text{KernelFunction}(|x - c_i| / KW)$$

Example



$$y^{est} = 2\phi_1(x) + 0.05\phi_2(x) + 0.5\phi_3(x)$$

where

$$\phi_i(x) = \text{KernelFunction}(|x - c_i| / KW)$$

RBFs with Linear Regression

All c_i 's are held constant (initialized randomly or on a grid in m-dimensional input space)

KW also held constant (initialized to be large enough that there's decent overlap between basis functions*)

*Usually much better than the crappy overlap on my diagram

$$y^{est} = 2\phi_1(x) + 0.05\phi_2(x) + 0.5\phi_3(x)$$

where

$$\phi_i(x) = \text{KernelFunction}(|x - c_i| / KW)$$

RBFs with Linear Regression

All c_i 's are held constant (initialized randomly or on a grid in m-dimensional input space)

KW also held constant (initialized to be large enough that there's decent overlap between basis functions*)

*Usually much better than the crappy overlap on my diagram

$$y^{est} = 2\phi_1(x) + 0.05\phi_2(x) + 5\phi_3(x)$$

where

$$\phi_i(x) = \text{KernelFunction}(|x - c_i| / KW)$$

then given Q basis functions, define the matrix Z such that $Z_{kj} = \text{KernelFunction}(|x_k - c_j| / KW)$ where x_k is the k th vector of inputs

And as before, $\beta = (Z^T Z)^{-1} (Z^T y)$

RBFs with NonLinear Regression

Allow the c_i 's to adapt to the data (initialized randomly or on a grid in m-dimensional input space)

KW allowed to adapt to the data. (Some folks even let each basis function have its own KW_j , permitting fine detail in dense regions of input space)

$$y^{est} = 2\phi_1(x) + 0.05\phi_2(x) + 5\phi_3(x)$$

where

$$\phi_i(x) = \text{KernelFunction}(|x - c_i| / KW)$$

But how do we now find all the β_j 's, c_i 's and KW ?

RBFs with NonLinear Regression

Allow the c_i 's to adapt to the data (initialized randomly or on a grid in m-dimensional input space)

KW allowed to adapt to the data. (Some folks even let each basis function have its own KW , permitting fine detail in dense regions of input space)

$$y^{est} = 2\phi_1(x) + 0.05\phi_2(x) + 5\phi_3(x)$$

where

$$\phi_i(x) = \text{KernelFunction}(|x - c_i| / KW)$$

But how do we now find all the β_i 's, c_i 's and KW ?

Answer: Gradient Descent

RBFs with NonLinear Regression

Allow the c_i 's to adapt to the data (initialized randomly or on a grid in m-dimensional input space)

KW allowed to adapt to the data. (Some folks even let each basis function have its own KW , permitting fine detail in dense regions of input space)

$$y^{est} = 2\phi_1(x) + 0.05\phi_2(x) + 5\phi_3(x)$$

where

$$\phi_i(x) = \text{KernelFunction}(|x - c_i| / KW)$$

But how do we now find all the β_i 's, c_i 's and KW ?

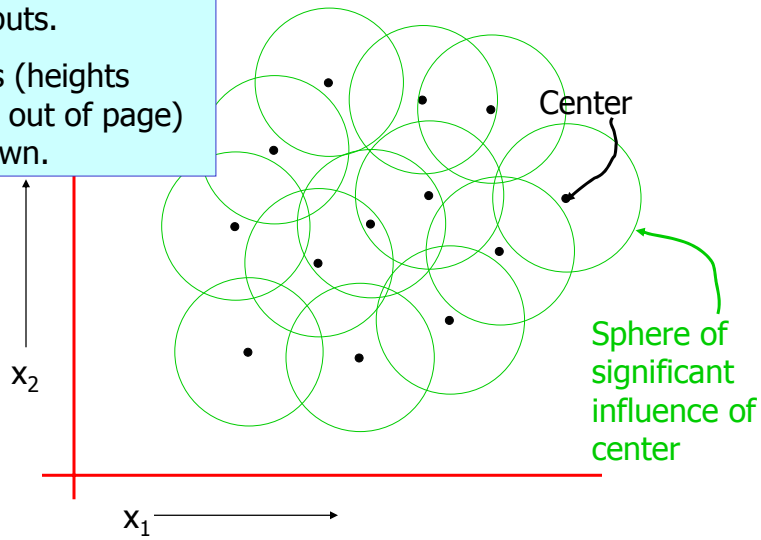
(But I'd like to see, or hope someone's already done, a hybrid, where the c_i 's and KW are updated with gradient descent while the β_i 's use matrix inversion)

Answer: Gradient Descent

Radial Basis Functions in 2-d

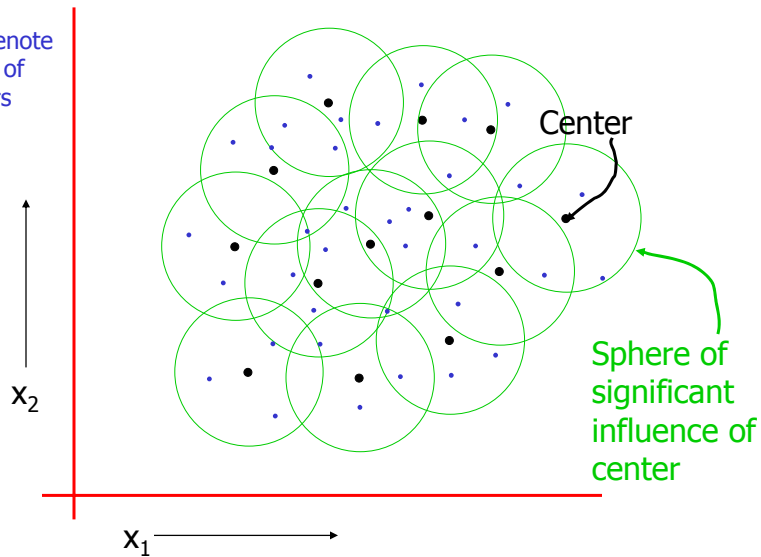
Two inputs.

Outputs (heights sticking out of page) not shown.



Happy RBFs in 2-d

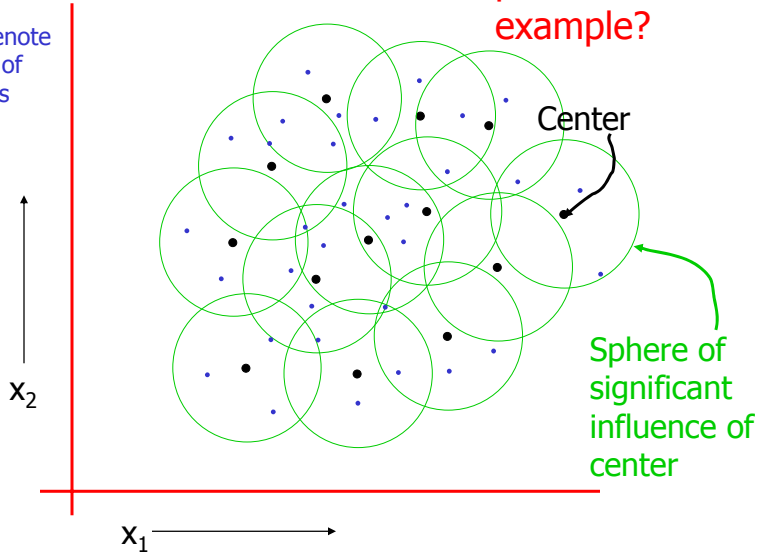
Blue dots denote coordinates of input vectors



Crabby RBFs in 2-d

What's the problem in this example?

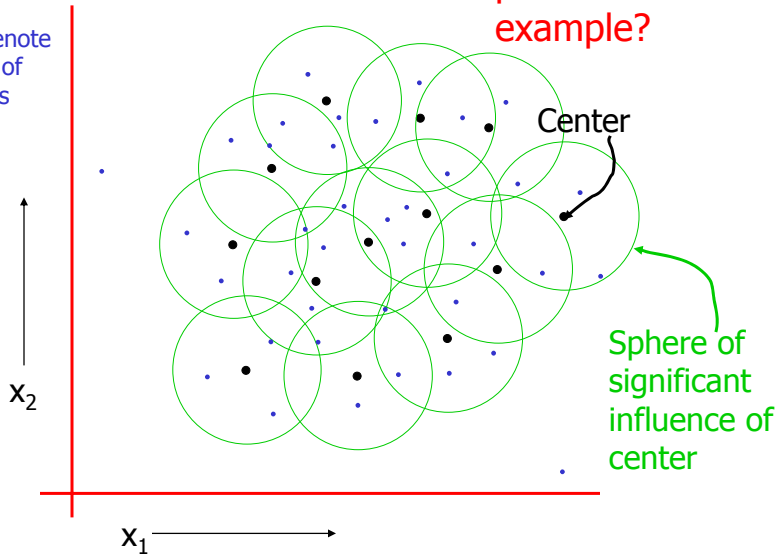
Blue dots denote coordinates of input vectors



More crabby RBFs

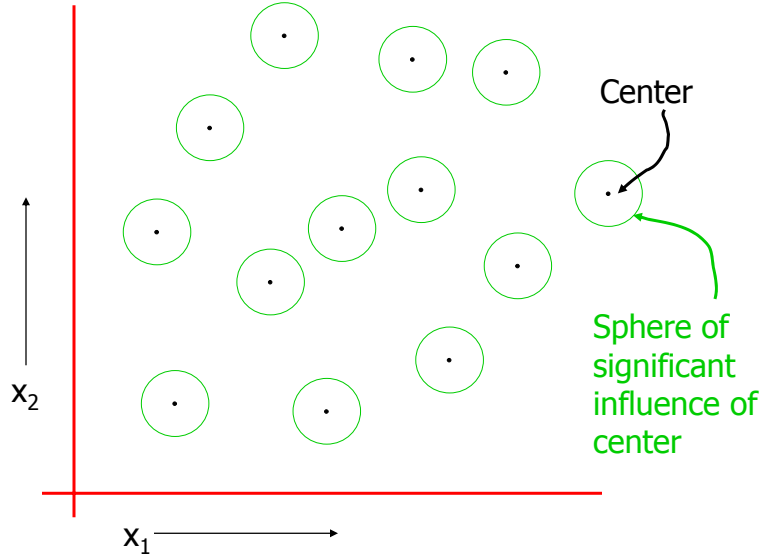
And what's the problem in this example?

Blue dots denote coordinates of input vectors



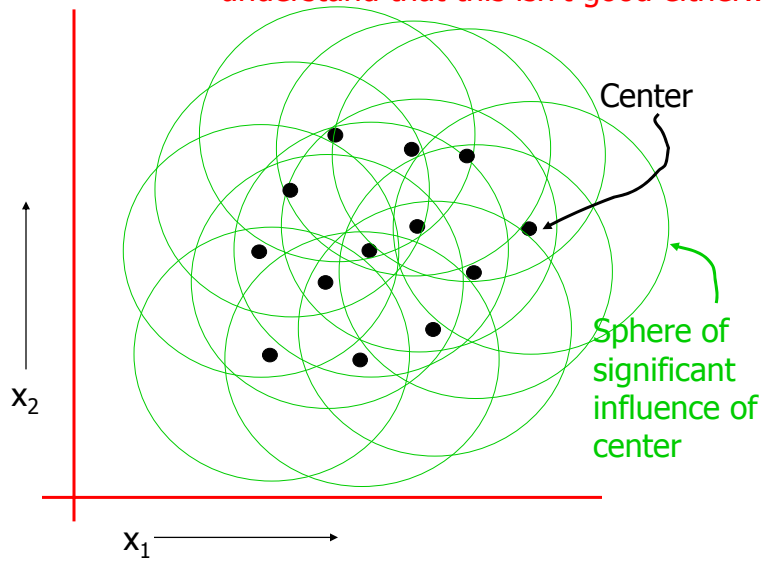
Hopeless!

Even before seeing the data, you should understand that this is a disaster!



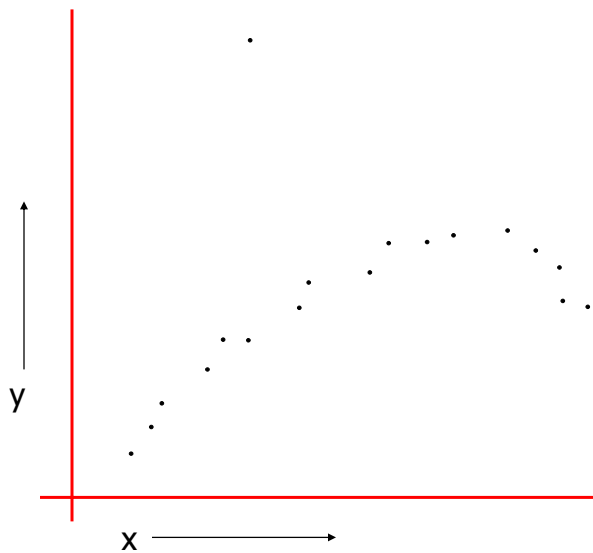
Unhappy

Even before seeing the data, you should understand that this isn't good either..



Robust Regression

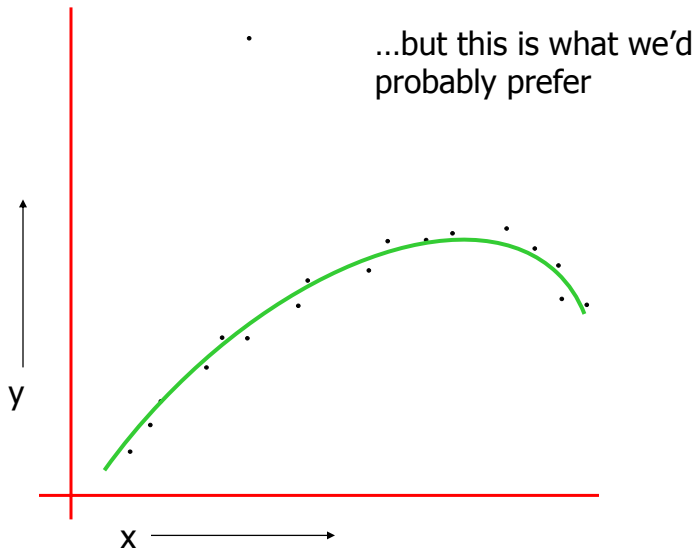
Robust Regression



Robust Regression

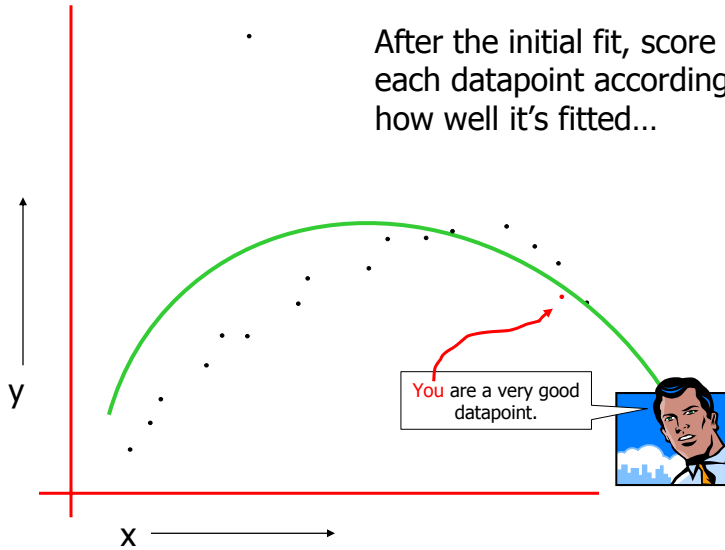


Robust Regression



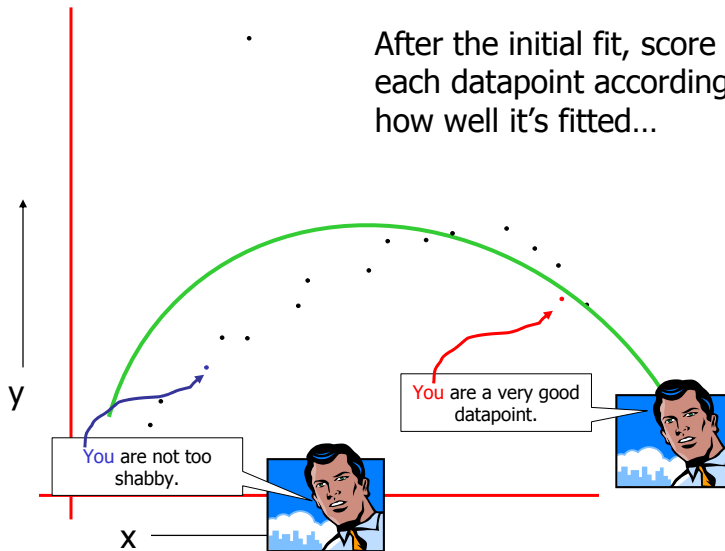
LOESS-based Robust Regression

After the initial fit, score each datapoint according to how well it's fitted...

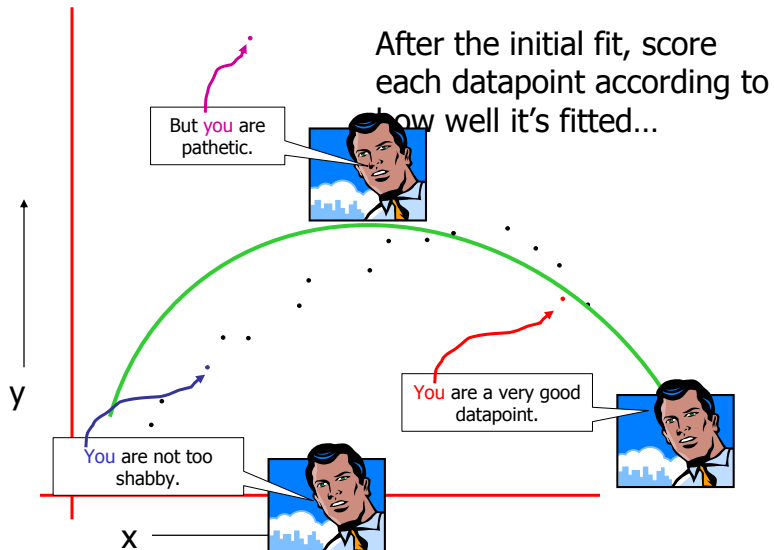


LOESS-based Robust Regression

After the initial fit, score each datapoint according to how well it's fitted...



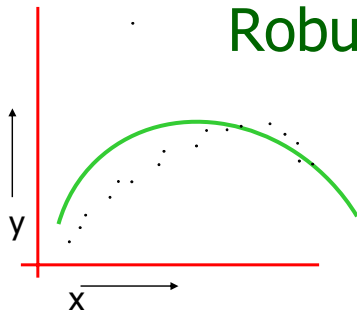
LOESS-based Robust Regression



Copyright © 2001, 2003, Andrew W. Moore

55

Robust Regression



For $k = 1$ to R ...

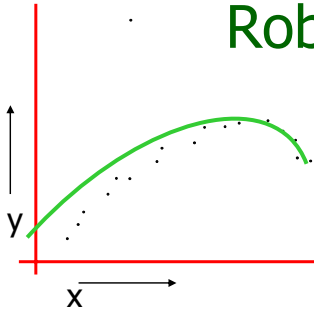
- Let (x_k, y_k) be the k th datapoint
- Let y_k^{est} be predicted value of y_k
- Let w_k be a weight for datapoint k that is large if the datapoint fits well and small if it fits badly:

$$w_k = \text{KernelFn}(|y_k - y_k^{est}|^2)$$

Copyright © 2001, 2003, Andrew W. Moore

56

Robust Regression



For $k = 1$ to $R...$

- Let (x_k, y_k) be the k th datapoint
- Let y_k^{est} be predicted value of y_k
- Let w_k be a weight for datapoint k that is large if the datapoint fits well and small if it fits badly:

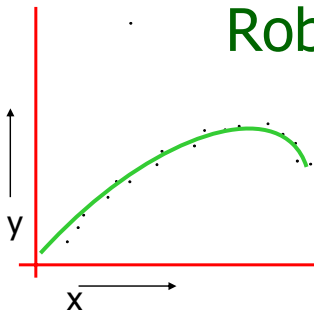
$$w_k = \text{KernelFn}([y_k - y_k^{est}]^2)$$

Then redo the regression using weighted datapoints.

Weighted regression was described earlier in the "vary noise" section, and is also discussed in the "Memory-based Learning" Lecture.

Guess what happens next?

Robust Regression



For $k = 1$ to $R...$

- Let (x_k, y_k) be the k th datapoint
- Let y_k^{est} be predicted value of y_k
- Let w_k be a weight for datapoint k that is large if the datapoint fits well and small if it fits badly:

$$w_k = \text{KernelFn}([y_k - y_k^{est}]^2)$$

Then redo the regression using weighted datapoints.

I taught you how to do this in the "Instance-based" lecture (only then the weights depended on distance in input-space)

Repeat whole thing until converged!

Robust Regression---what we're doing

What regular regression does:

Assume y_k was originally generated using the following recipe:

$$y_k = \beta_0 + \beta_1 x_k + \beta_2 x_k^2 + N(0, \sigma^2)$$

Computational task is to find the Maximum Likelihood β_0 , β_1 and β_2

Robust Regression---what we're doing

What LOESS robust regression does:

Assume y_k was originally generated using the following recipe:

With probability p :

$$y_k = \beta_0 + \beta_1 x_k + \beta_2 x_k^2 + N(0, \sigma^2)$$

But otherwise

$$y_k \sim N(\mu, \sigma_{huge}^2)$$

Computational task is to find the Maximum Likelihood β_0 , β_1 , β_2 , p , μ and σ_{huge}

Robust Regression---what we're doing

What LOESS robust regression does:

Assume y_k was originally generated using the following recipe:

With probability p :

$$y_k = \beta_0 + \beta_1 x_k + \beta_2 x_k^2 + N(0, \sigma^2)$$

But otherwise

$$y_k \sim N(\mu, \sigma_{huge}^2)$$

Computational task is to find the Maximum Likelihood $\beta_0, \beta_1, \beta_2, p, \mu$ and σ_{huge}

Mysteriously, the reweighting procedure does this computation for us.

Your first glimpse of two spectacular letters:

E.M.

Regression Trees

Regression Trees

- "Decision trees for regression"

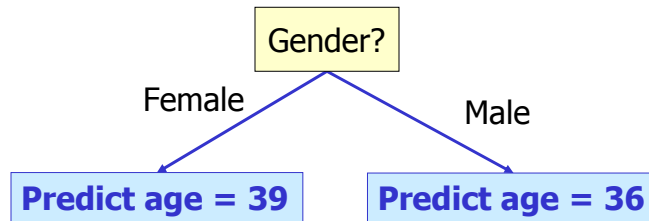
A regression tree leaf



Predict age = 47

Mean age of records
matching this leaf node

A one-split regression tree



Choosing the attribute to split on

Gender	Rich?	Num. Children	Num. Beany Babies	Age
Female	No	2	1	38
Male	No	0	0	24
Male	Yes	0	5+	72
:	:	:	:	:

- We can't use information gain.
- What should we use?

Choosing the attribute to split on

Gender	Rich?	Num. Children	Num. Beany Babies	Age
Female	No	2	1	38
Male	No	0	0	24
Male	Yes	0	5+	72
:	:	:	:	:

$MSE(Y|X)$ = The expected squared error if we must predict a record's Y value given only knowledge of the record's X value

If we're told $x=j$, the smallest expected error comes from predicting the mean of the Y-values among those records in which $x=j$. Call this mean quantity $\mu_y^{x=j}$

Then...

$$MSE(Y | X) = \frac{1}{R} \sum_{j=1}^{N_x} \sum_{(k \text{ such that } x_k=j)} (y_k - \mu_y^{x=j})^2$$

Choosing the attribute to split on

Gender	Rich?	Num. Children	Num. Beany Babies	Age
Female	N			
Male	N			
Male	Y			
:	:			

Regression tree attribute selection: greedily choose the attribute that minimizes $MSE(Y|X)$

Guess what we do about real-valued inputs?

Guess how we prevent overfitting

$MSE(Y|X)$ = The expected squared error if we must predict a record's Y value given only knowledge of the record's X value

If we're told $x=j$, the smallest expected error comes from predicting the mean of the Y-values among those records in which $x=j$. Call this mean quantity $\mu_y^{x=j}$

Then...

$$MSE(Y | X) = \frac{1}{R} \sum_{j=1}^{N_x} \sum_{(k \text{ such that } x_k=j)} (y_k - \mu_y^{x=j})^2$$

Pruning Decision

...property-owner = Yes

Gender?

Do I deserve to live?

Female

Male

Predict age = 39

Predict age = 36

property-owning females = 56712

Mean age among POFs = 39

Age std dev among POFs = 12

property-owning males = 55800

Mean age among POMs = 36

Age std dev among POMs = 11.5

Use a standard Chi-squared test of the null-hypothesis "these two populations have the same mean" and Bob's your uncle.

Linear Regression Trees

...property-owner = Yes

Gender?

Also known as "Model Trees"

Female

Male

Predict age =
 $26 + 6 * \text{NumChildren} - 2 * \text{YearsEducation}$

Predict age =
 $24 + 7 * \text{NumChildren} - 2.5 * \text{YearsEducation}$

Leaves contain linear functions (trained using linear regression on all records matching that leaf)

Split attribute chosen to minimize MSE of regressed children.

Pruning with a different Chi-squared

Linear Regression Trees

...property-owner = Yes

Gender?

Female

Predict age =

$$26 + 6 * M + 2 * Year$$

Also known as
"M... Trees"

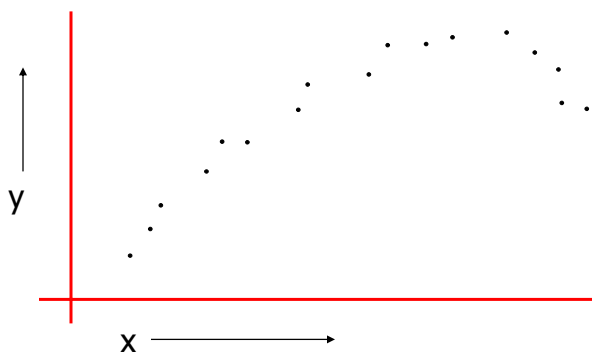
Detail: You typically ignore any categorical attribute that has been tested on higher up in the tree during the regression. But use all untested attributes, and use real-valued attributes even if they've been tested above

Leaves contain linear regression functions (trained on records matching the leaf) chosen to minimize regression error.

Pruning with a different Chi-squared

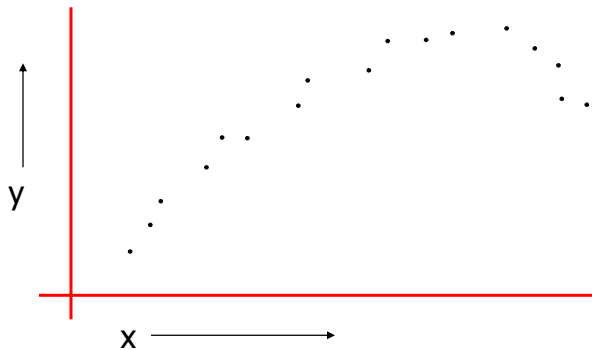
Test your understanding

Assuming **regular** regression trees, can you sketch a graph of the fitted function $y^{est}(x)$ over this diagram?



Test your understanding

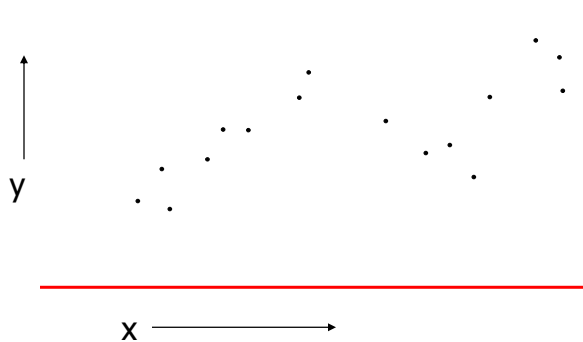
Assuming **linear** regression trees, can you sketch a graph of the fitted function $y^{est}(x)$ over this diagram?



Multilinear Interpolation

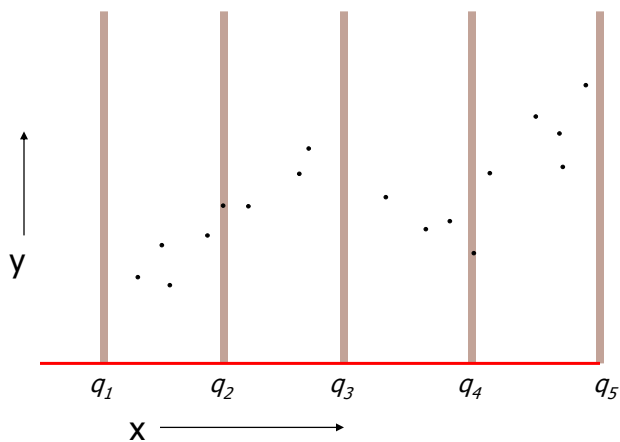
Multilinear Interpolation

Consider this dataset. Suppose we wanted to create a continuous and piecewise linear fit to the data



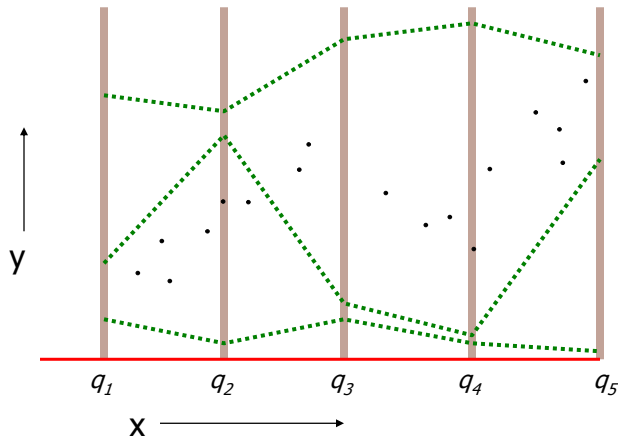
Multilinear Interpolation

Create a set of knot points: selected X-coordinates (usually equally spaced) that cover the data



Multilinear Interpolation

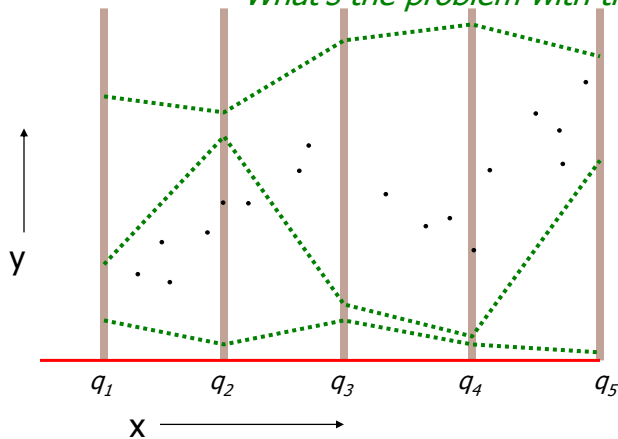
We are going to assume the data was generated by a noisy version of a function that can only bend at the knots. Here are 3 examples (none fits the data well)



How to find the best fit?

Idea 1: Simply perform a separate regression in each segment for each part of the curve

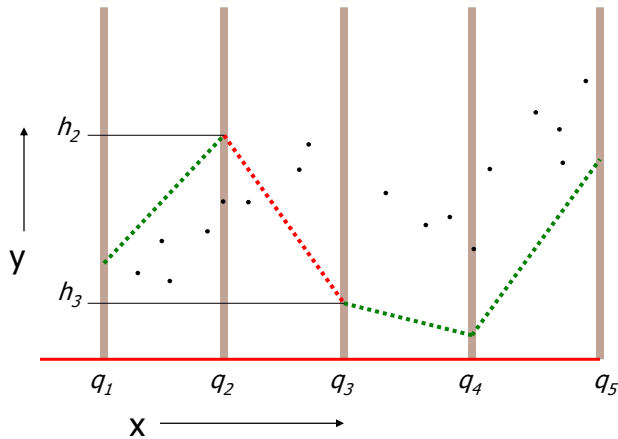
What's the problem with this idea?



How to find the best fit?

Let's look at what goes on in the red segment

$$y^{est}(x) = \frac{(q_3 - x)}{w} h_2 + \frac{(q_2 - x)}{w} h_3 \text{ where } w = q_3 - q_2$$

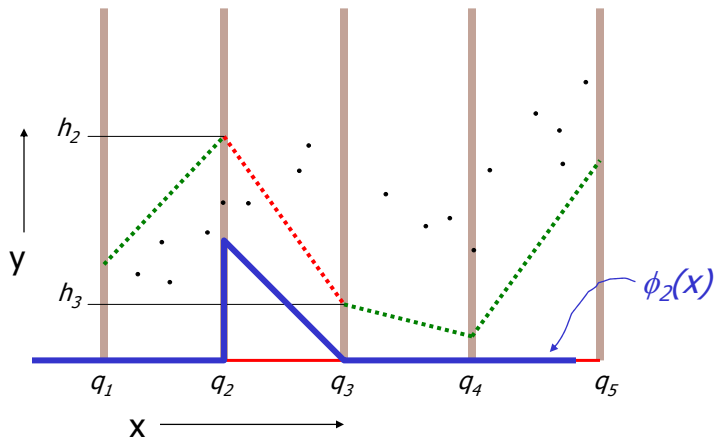


How to find the best fit?

In the red segment...

$$y^{est}(x) = h_2 \phi_2(x) + h_3 \phi_3(x)$$

$$\text{where } \phi_2(x) = 1 - \frac{x - q_2}{w}, \phi_3(x) = 1 - \frac{q_3 - x}{w}$$

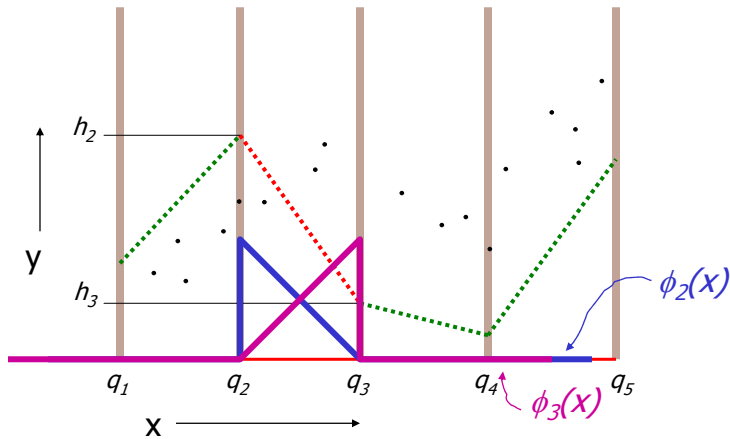


How to find the best fit?

In the red segment...

$$y^{est}(x) = h_2\phi_2(x) + h_3\phi_3(x)$$

$$\text{where } \phi_2(x) = 1 - \frac{x - q_2}{w}, \phi_3(x) = 1 - \frac{q_3 - x}{w}$$

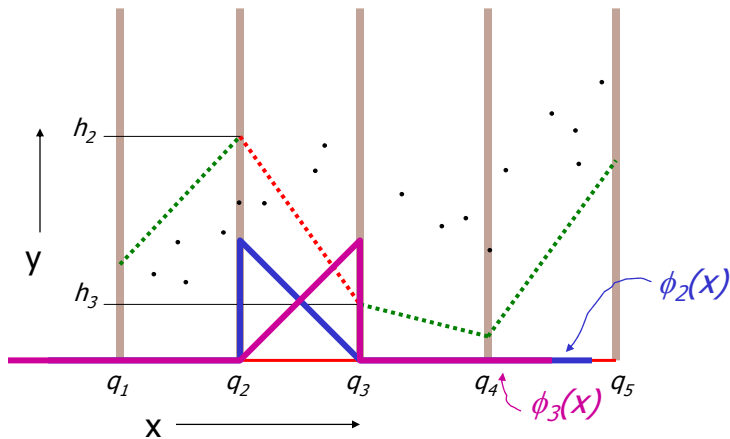


How to find the best fit?

In the red segment...

$$y^{est}(x) = h_2\phi_2(x) + h_3\phi_3(x)$$

$$\text{where } \phi_2(x) = 1 - \frac{|x - q_2|}{w}, \phi_3(x) = 1 - \frac{|x - q_3|}{w}$$

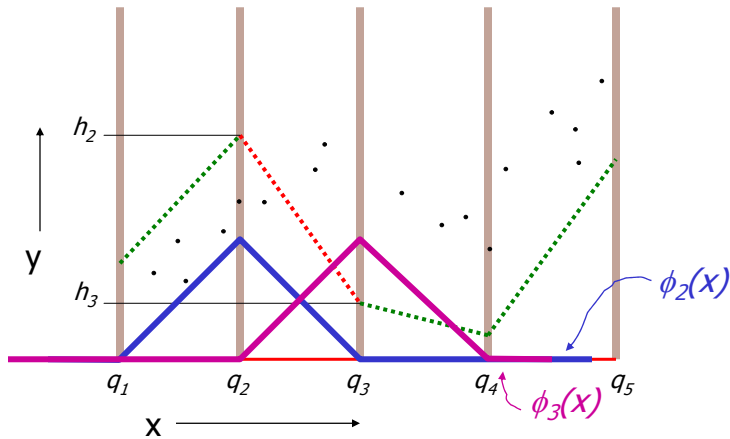


How to find the best fit?

In the red segment...

$$y^{est}(x) = h_2\phi_2(x) + h_3\phi_3(x)$$

$$\text{where } \phi_2(x) = 1 - \frac{|x - q_2|}{w}, \phi_3(x) = 1 - \frac{|x - q_3|}{w}$$

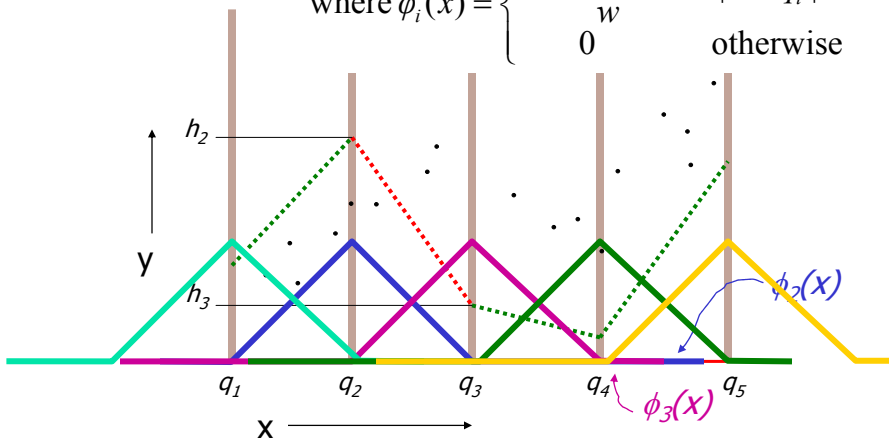


How to find the best fit?

In **general**

$$y^{est}(x) = \sum_{i=1}^{N_K} h_i\phi_i(x)$$

$$\text{where } \phi_i(x) = \begin{cases} 1 - \frac{|x - q_i|}{w} & \text{if } |x - q_i| < w \\ 0 & \text{otherwise} \end{cases}$$

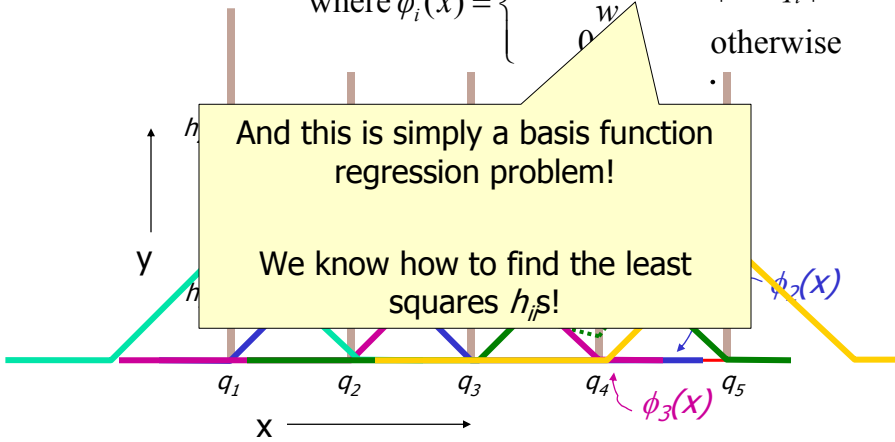


How to find the best fit?

In **general**

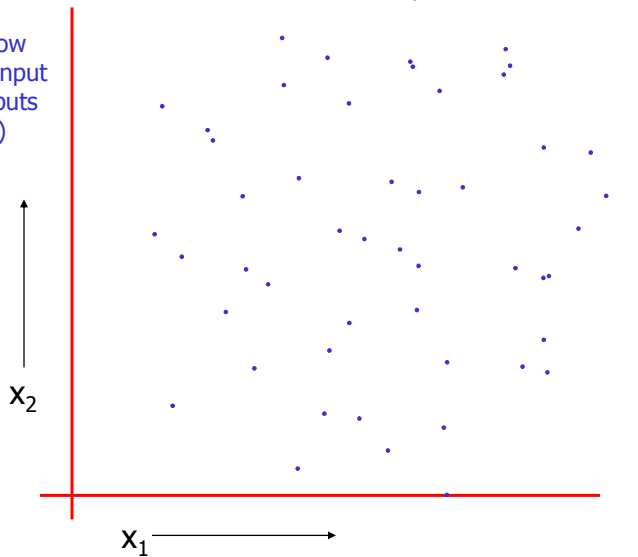
$$y^{est}(x) = \sum_{i=1}^{N_k} h_i \phi_i(x)$$

$$\text{where } \phi_i(x) = \begin{cases} 1 - \frac{|x - q_i|}{w} & \text{if } |x - q_i| < w \\ 0 & \text{otherwise} \end{cases}$$



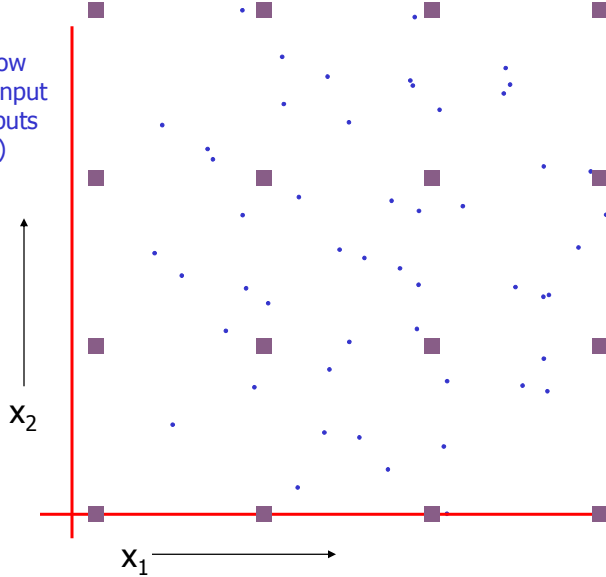
In two dimensions...

Blue dots show locations of input vectors (outputs not depicted)



In two dimensions...

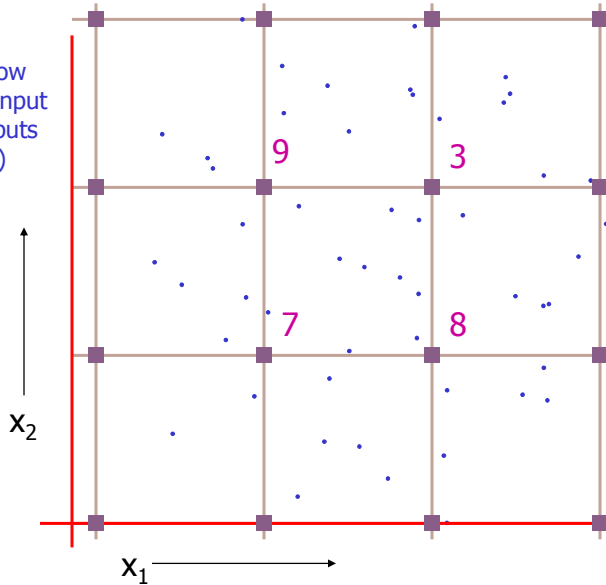
Blue dots show locations of input vectors (outputs not depicted)



Each purple dot is a knot point. It will contain the height of the estimated surface

In two dimensions...

Blue dots show locations of input vectors (outputs not depicted)



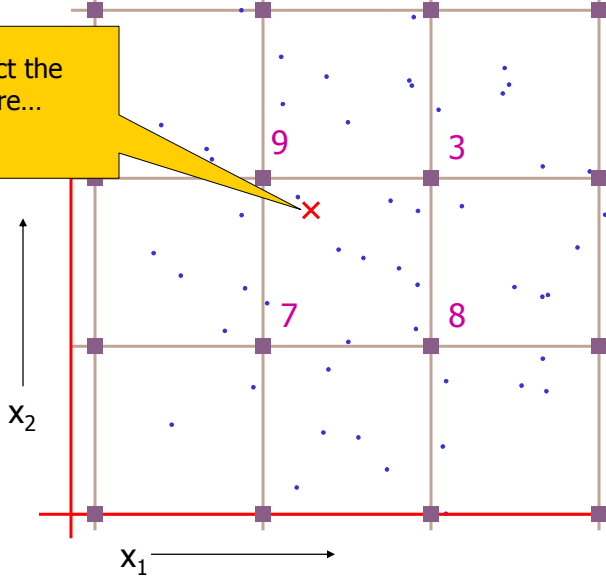
Each purple dot is a knot point. It will contain the height of the estimated surface

But how do we do the interpolation to ensure that the surface is continuous?

In two dimensions...

Bl
lo
ve
nc

To predict the value here...



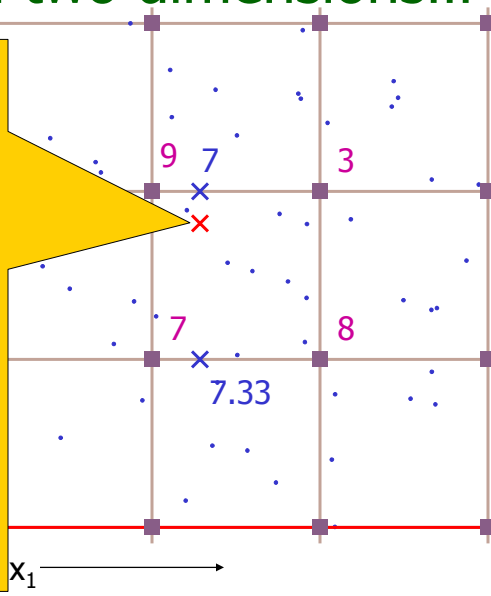
Each purple dot is a knot point. It will contain the height of the estimated surface

But how do we do the interpolation to ensure that the surface is continuous?

In two dimensions...

Bl
lo
ve
nc

To predict the value here...
First interpolate its value on two opposite edges...



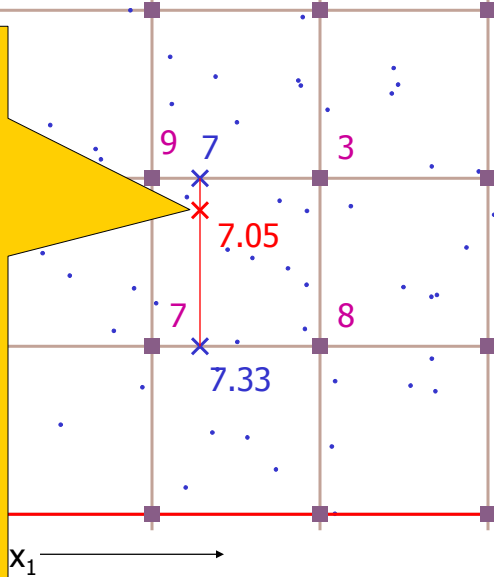
Each purple dot is a knot point. It will contain the height of the estimated surface

But how do we do the interpolation to ensure that the surface is continuous?

In two dimensions...

Bl
lo
ve
nc

To predict the value here...
First interpolate its value on two opposite edges...
Then interpolate between those two values



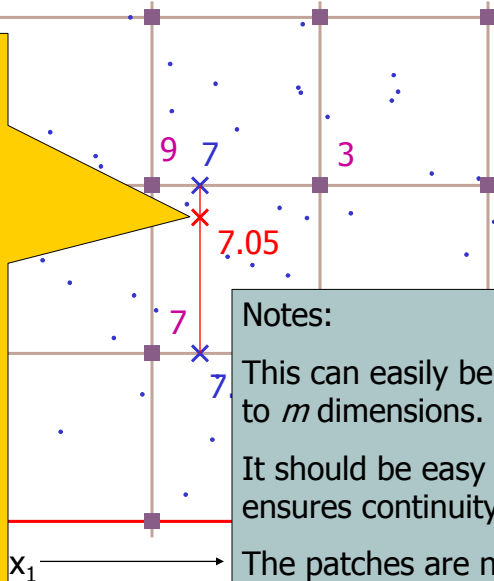
Each purple dot is a knot point. It will contain the height of the estimated surface

But how do we do the interpolation to ensure that the surface is continuous?

In two dimensions...

Bl
lo
ve
nc

To predict the value here...
First interpolate its value on two opposite edges...
Then interpolate between those two values



Each purple dot is a knot point. It will contain the height of the estimated surface

But how do we do the interpolation to ensure that the

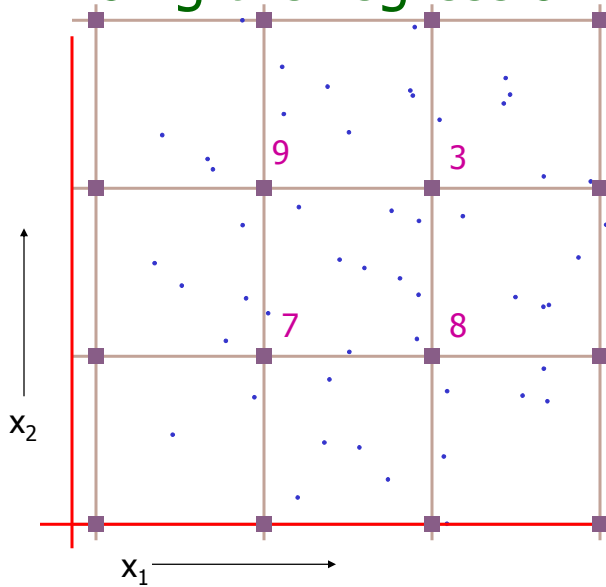
Notes:

This can easily be generalized to m dimensions.

It should be easy to see that it ensures continuity

The patches are not linear

Doing the regression



Given data, how do we find the optimal knot heights?

Happily, it's simply a two-dimensional basis function problem.

(Working out the basis functions is tedious, unilluminating, and easy)

What's the problem in higher dimensions?

MARS: Multivariate Adaptive Regression Splines

MARS

- Multivariate Adaptive Regression Splines
- Invented by Jerry Friedman (one of Andrew's heroes)
- Simplest version:

Let's assume the function we are learning is of the following form:

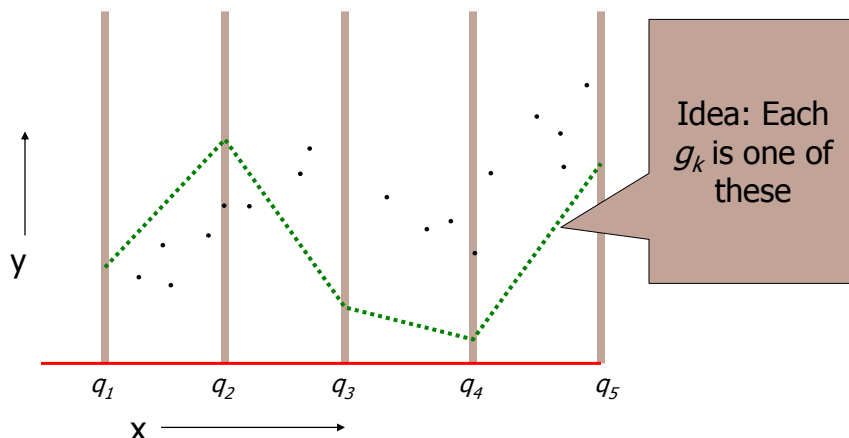
$$y^{est}(\mathbf{x}) = \sum_{k=1}^m g_k(x_k)$$

Instead of a linear combination of the inputs, it's a linear combination of non-linear functions of *individual* inputs

MARS

$$y^{est}(\mathbf{x}) = \sum_{k=1}^m g_k(x_k)$$

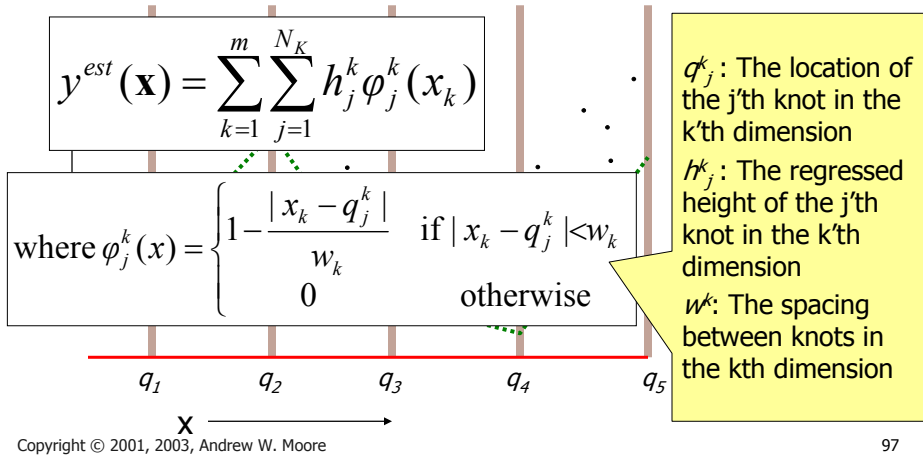
Instead of a linear combination of the inputs, it's a linear combination of non-linear functions of *individual* inputs



MARS

$$y^{est}(\mathbf{x}) = \sum_{k=1}^m g_k(x_k)$$

Instead of a linear combination of the inputs, it's a linear combination of non-linear functions of *individual* inputs



That's not complicated enough!

- Okay, now let's get serious. We'll allow arbitrary "two-way interactions":

$$y^{est}(\mathbf{x}) = \sum_{k=1}^m g_k(x_k) + \sum_{k=1}^m \sum_{t=k+1}^m g_{kt}(x_k, x_t)$$

The function we're learning is allowed to be a sum of non-linear functions over all one-d and 2-d subsets of attributes

Can still be expressed as a linear combination of basis functions

Thus learnable by linear regression

Full MARS: Uses cross-validation to choose a subset of subspaces, knot resolution and other parameters.

If you like MARS...

...See also CMAC (Cerebellar Model Articulated Controller) by James Albus (another of Andrew's heroes)

- Many of the same gut-level intuitions
- But entirely in a neural-network, biologically plausible way
 - (All the low dimensional functions are by means of lookup tables, trained with a delta-rule and using a clever blurred update and hash-tables)

Where are we now?

Inputs	Inference Engine Learn	$p(E_1 E_2)$	Joint DE
Inputs	Classifier	Predict category	Dec Tree, Gauss/Joint BC, Gauss Naïve BC,
Inputs	Density Estimator	Probability	Joint DE, Naïve DE, Gauss/Joint DE, Gauss Naïve DE
Inputs	Regressor	Predict real no.	Linear Regression, Polynomial Regression, RBFs, Robust Regression Regression Trees, Multilinear Interp, MARS

Citations

Radial Basis Functions

T. Poggio and F. Girosi,
Regularization Algorithms for
Learning That Are Equivalent to
Multilayer Networks, *Science*,
247, 978--982, 1989

LOESS

W. S. Cleveland, Robust Locally
Weighted Regression and
Smoothing Scatterplots, *Journal
of the American Statistical
Association*, 74, 368, 829-836,
December, 1979

Regression Trees etc

L. Breiman and J. H. Friedman and
R. A. Olshen and C. J. Stone,
Classification and Regression
Trees, Wadsworth, 1984

J. R. Quinlan, Combining Instance-
Based and Model-Based
Learning, *Machine Learning:
Proceedings of the Tenth
International Conference*, 1993

MARS

J. H. Friedman, Multivariate
Adaptive Regression Splines,
Department for Statistics,
Stanford University, 1988,
Technical Report No. 102