# CPSC 340:
# Machine Learning and Data Mining

Non-Linear Regression

Fall 2016

# Admin

- Assignment 2 is due now.
  - 1 late day to hand it in on Wednesday, 2 for Friday, 3 for next Monday.
- Assignment 3 will be out by early next week.
  - Due October 19 (so we can release solutions before the midterm).
- We will have tutorials on Tuesday/Wednesday of next week:
  - Focusing on multivariate calculus in matrix notation.
- Tutorial room change: T1D (Monday @5pm) moved to DMP 101.

THE DATA BEHIND MASSIVE
OPEN ONLINE COURSES (MOOCS) AT
U D A C I T Y
Thursday October 13, 2016 at 5:30pm

.ıIDataSense   BIG DATA UNIVERSITY

*"In this talk, Eli will provide an overview of how big data is used to create and power
Silicon Valley's greatest companies, with specific examples from Udacity."*

Location TBA!

There will be free food and drinks. More info and to get your tickets:
https://goo.gl/QGsnUU

https://www.facebook.com/events/645894388926612/

# Last Time: Linear Regression

- We discussed linear models:

$$y_i = w_1 x_{i1} + w_2 x_{i2} + \cdots + w_d x_{id}$$

$$= \sum_{j=1}^{d} w_j x_{ij} = w^T x_i$$

- "Multiply feature $x_{ij}$ by weight $w_j$, add them to get $y_i$".
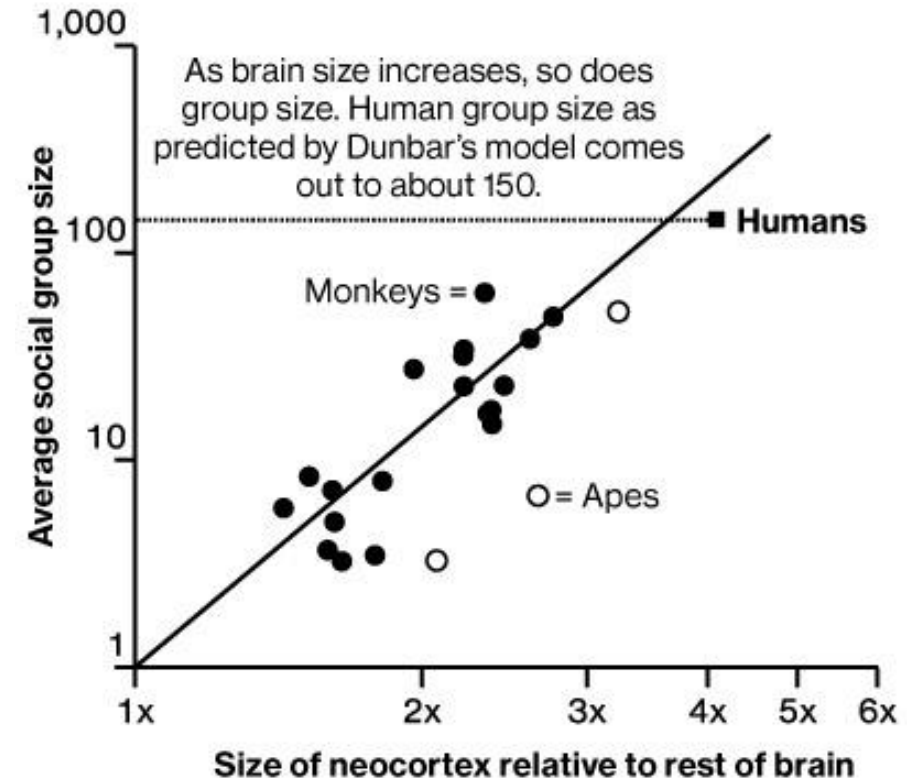
- We discussed squared error function:

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2$$

Predicted value ← → True value

- Interactive demo:
  - http://setosa.io/ev/ordinary-least-squares-regression



**The Social Cortex**

1,000

As brain size increases, so does group size. Human group size as predicted by Dunbar's model comes out to about 150.

100 ······························ ■ **Humans**

Average social group size

Monkeys = ●

10

O = Apes

1

1x    2x    3x    4x    5x    6x

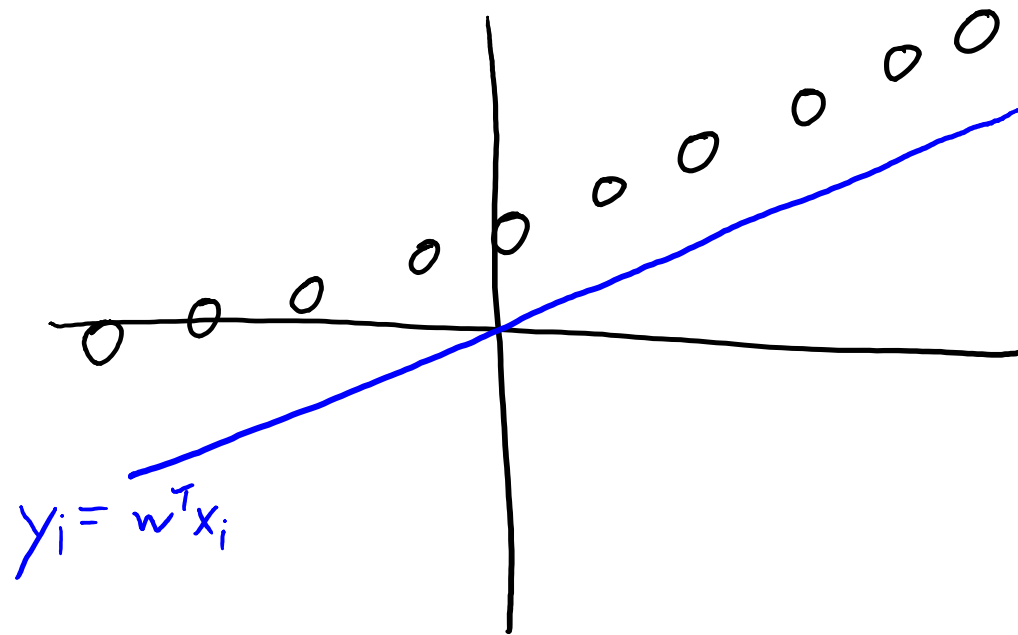**Size of neocortex relative to rest of brain**

DATA: THE SOCIAL BRAIN HYPOTHESIS, DUNBAR 1998

To predict on test case $\hat{x}_i$ use $\hat{y}_i = w^T \hat{x}_i$
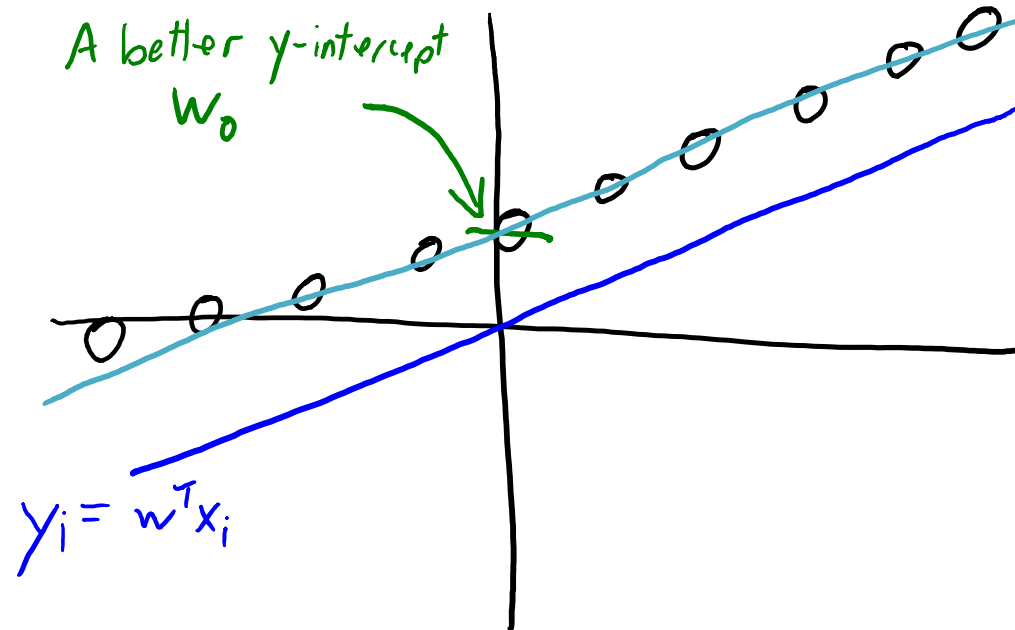
# Why don't we have a y-intercept?

- Last time: Linear models with no y-intercept.
  - Linear model is $y_i = w^T x_i$ instead of $y_i = w^T x_i + w_0$ with y-intercept $w_0$.
  - So if $x_i = 0$ then we must predict $y_i = 0$.

# Why don't we have a y-intercept?

- Last time: Linear models with no y-intercept.
  - Linear model is $y_i = w^T x_i$ instead of $y_i = w^T x_i + w_0$ with y-intercept $w_0$.
  - So if $x_i = 0$ then we must predict $y_i = 0$.

Adding
y-intercept
fixes this.

A better y-intercept
$w_0$

$y_i = w^T x_i + w_0$

Even "least squares"
solution must go
through origin.

$y_i = w^T x_i$

# Adding a Bias Variable

- Simple trick to add a y-intercept ("bias") variable:
  - Make a new matrix "Z" with an extra feature that is always "1".

$$X = \begin{bmatrix} 0.1 & 0.3 \\ 0.5 & -0.6 \\ 0.2 & 0.4 \end{bmatrix} \qquad Z = \begin{bmatrix} 1 & 0.1 & 0.3 \\ 1 & 0.5 & -0.6 \\ 1 & 0.2 & 0.4 \end{bmatrix}$$

$$X$$

- Now use "Z" as features to get a model with a non-zero y-intercept:

$$y_i = w_0 \, z_{i0} + w_1 \, z_{i1} + w_2 \, z_{i2}$$

$$\hookrightarrow "1" \qquad \hookrightarrow x_{i1} \qquad \hookrightarrow x_{i2}$$

$$= w_0 + w_1 \, x_{i1} + w_2 \, x_{i2}$$

- So we can have a non-zero y-intercept by changing features.

# Linear Least Squares

Training:
$$w = (X' * X) \setminus (X' * y)$$

Prediction:
$$y = X * w$$

Why?

$$y_i = w^T x_i \quad \text{so} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_n \end{bmatrix} = \begin{bmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{bmatrix} = \begin{bmatrix} \text{---} x_1^T \text{---} \\ \text{---} x_2^T \text{---} \\ \vdots \\ \text{---} x_n^T \text{---} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_d \end{bmatrix} = Xw$$

$X$

# Linear Least Squares

Training:

$$w = (X' * X) \backslash (X' * y)$$

Why?

Want 'w' that minimizes

$$f(w) = \frac{1}{2}\sum_{i=1}^{n}(w^Tx_i - y_i)^2 = \frac{1}{2}\sum_{i=1}^{n}r_i^2 = \frac{1}{2}r^Tr = \frac{1}{2}\|r\|_2^2 = \boxed{\frac{1}{2}\|Xw - y\|^2}$$

$r_i$

Define "residual" $r_i$ as signed error on example 'i':

$$r_i = w^Tx_i - y_i$$

$$r = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{bmatrix} = \begin{bmatrix} w^Tx_1 - y_1 \\ w^Tx_2 - y_2 \\ \vdots \\ w^Tx_n - y_n \end{bmatrix} = \underbrace{\begin{bmatrix} w^Tx_1 \\ w^Tx_2 \\ \vdots \\ w^Tx_n \end{bmatrix}}_{Xw} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \boxed{Xw - y}$$

# Linear Least Squares

$$\|r\|_2 = \sqrt{\sum_{j=1}^{t} r_j^2}$$

Training:

$$w = (X' * X) \setminus (X' * y)$$

Why?

$$r^T r = \|r\|^2$$

Want 'w' that <u>minimizes</u>

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2 = \frac{1}{2} \|Xw - y\|_2^2 = \frac{1}{2}(Xw - y)^T(Xw - y)$$

Let's expand then compute gradient.

$$= \frac{1}{2}\left((Xw)^T - y^T\right)(Xw - y)$$

$$(AB)^T = B^T A^T$$

$$= \frac{1}{2}\left(\underline{w^T X^T} - y^T\right)(Xw - y)$$

"distributive"

$$= \frac{1}{2}\left(w^T X^T(Xw - y) - y^T(Xw - y)\right)$$

A <u>good</u> way to <u>check</u> your step: make sure that dimensions all make sense.

$$= \frac{1}{2}\left(w^T X^T Xw - w^T X^T y - y^T Xw + y^T y\right)$$

$$= \frac{1}{2} w^T X^T Xw - w^T X^T y + \frac{1}{2} y^T y$$

$$w^T v = v^T w$$
$$v = X^T y$$

# Linear Least Squares

Training:

$$w = (X' * X) \setminus (X' * y)$$

Why?

Want 'w' that minimizes

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2 = \frac{1}{2} \| Xw - y \|_2^2 = \frac{1}{2} w^T X^T X w - w^T X^T y + \frac{1}{2} y^T y \left\{ \begin{array}{l} \text{a quadratic function} \\ \text{(in matrix notation)} \end{array} \right.$$

$A$

$ax^2 + bx + c$

What are the gradients of these terms?

$$\nabla f(w) = X^T X w - X^T y + 0$$

So at a minimizer where $\nabla f(w) = 0$

we have: $\boxed{X^T X w = X^T y}$

See note on webpage for derivation

Cheat sheet: $\nabla_w [c] = 0$

$\nabla_w [w^T b] = b$

$\nabla_w [\frac{1}{2} w^T A w] = Aw$  for symmetric $A$.

Like $\frac{d}{dx}[\frac{1}{2} ax^2]$

$= ax$

# Linear Least Squares

Training:

$$\overbrace{\phantom{(X'*X)}}^{A} \quad \underset{\uparrow}{\text{"solve linear system"}} \quad \overbrace{\phantom{(X'*y)}}^{b}$$

$$w = \underbrace{(X'*X) \diagdown (X'*y)}_{\text{Why?}}$$

Want 'w' that <u>minimizes</u>

$$f(w) = \frac{1}{2}\sum_{i=1}^{n}(w^T x_i - y_i)^2 = \frac{1}{2}\|Xw - y\|_2^2$$

We know that minimizer must have

$$\underbrace{X^T X}_{\text{Some matrix}} w = \underbrace{X^T y}_{\text{Some vector}} \qquad \text{"normal equations"}$$

This is a linear system $Aw = b$

for some matrix 'A' and vector 'b'

Note that
$f(w)$ is a "convex" function
so solving $\nabla f(w) = 0$ gives minimizer

# Least Squares Issues

- Issues with least squares model:
  - Solution might not be unique.
  - It is sensitive to outliers.
  - It always uses all features.
  - Data can might so big we can't store $X^TX$.
  - It might predict outside range of $y_i$ values.
  - It assumes a linear relationship between $x_i$ and $y_i$.

$X$ is $n \times d$
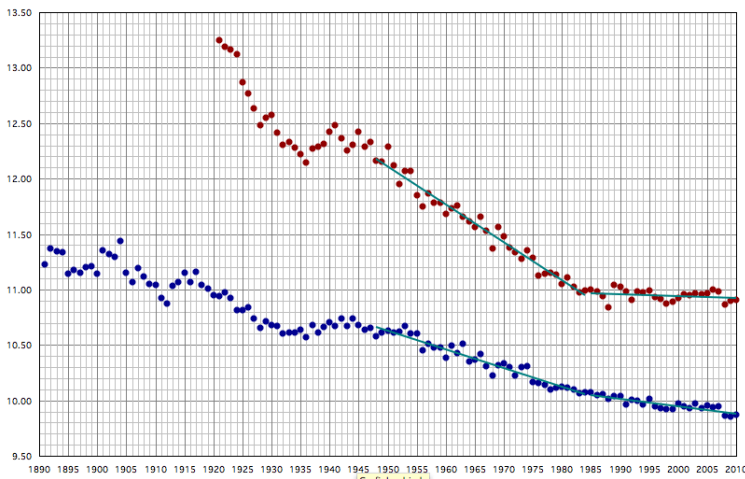
so $X^T$ is $d \times n$

and $X^TX$ is $d \times d$.

Costs $O(nd^2)$ to calculate:

- Each of the $O(d^2)$ elements is an inner product between length '$n$' vectors.
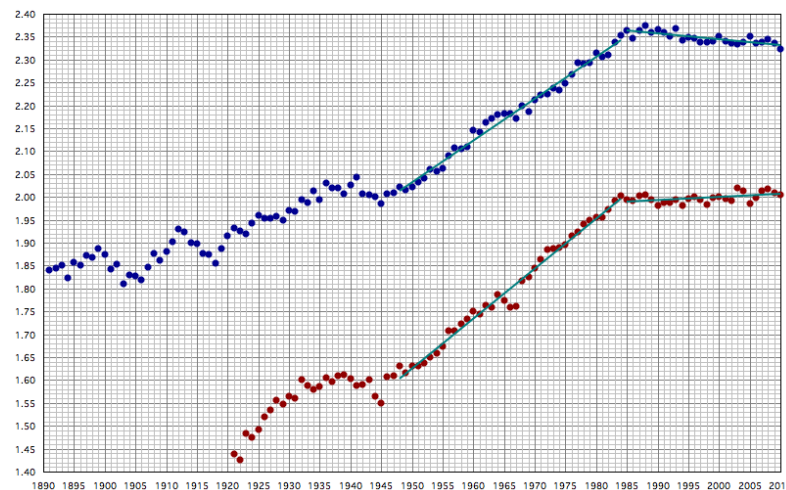
# Example: Non-Linear Progressions in Athletics

- Are top athletes going faster, higher, and farther?
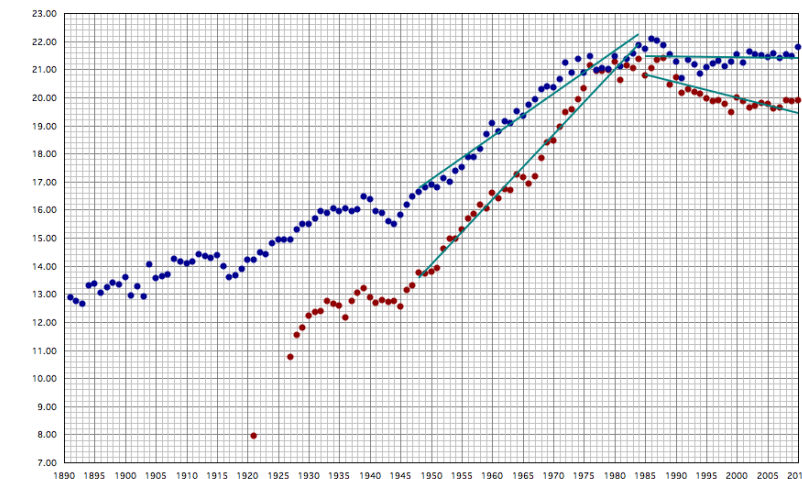


100m PROGRESSION MEN AND WOMEN (mean of top ten)



HIGH JUMP PROGRESSION MEN AND WOMEN (mean of top ten)



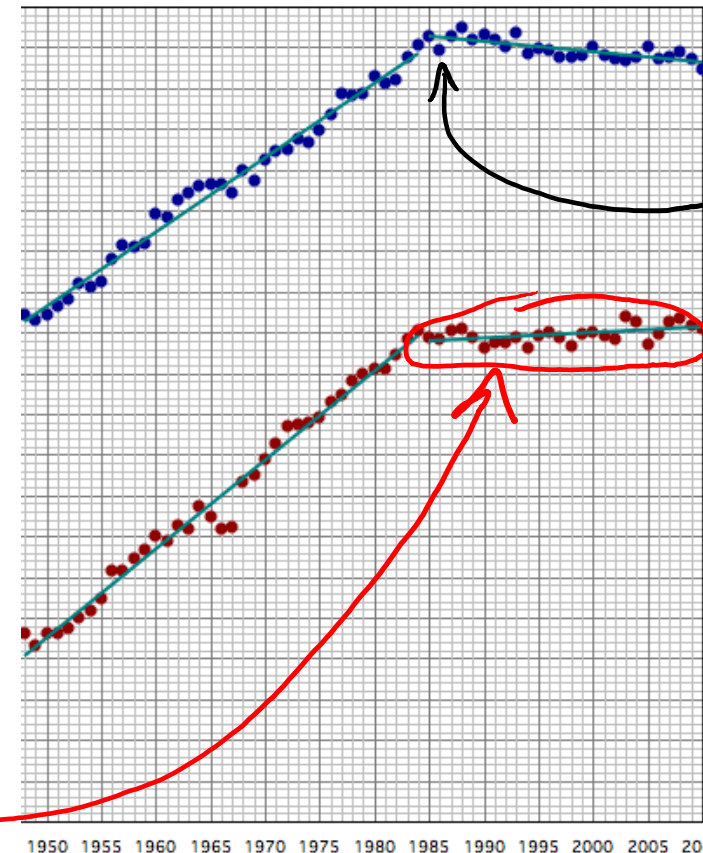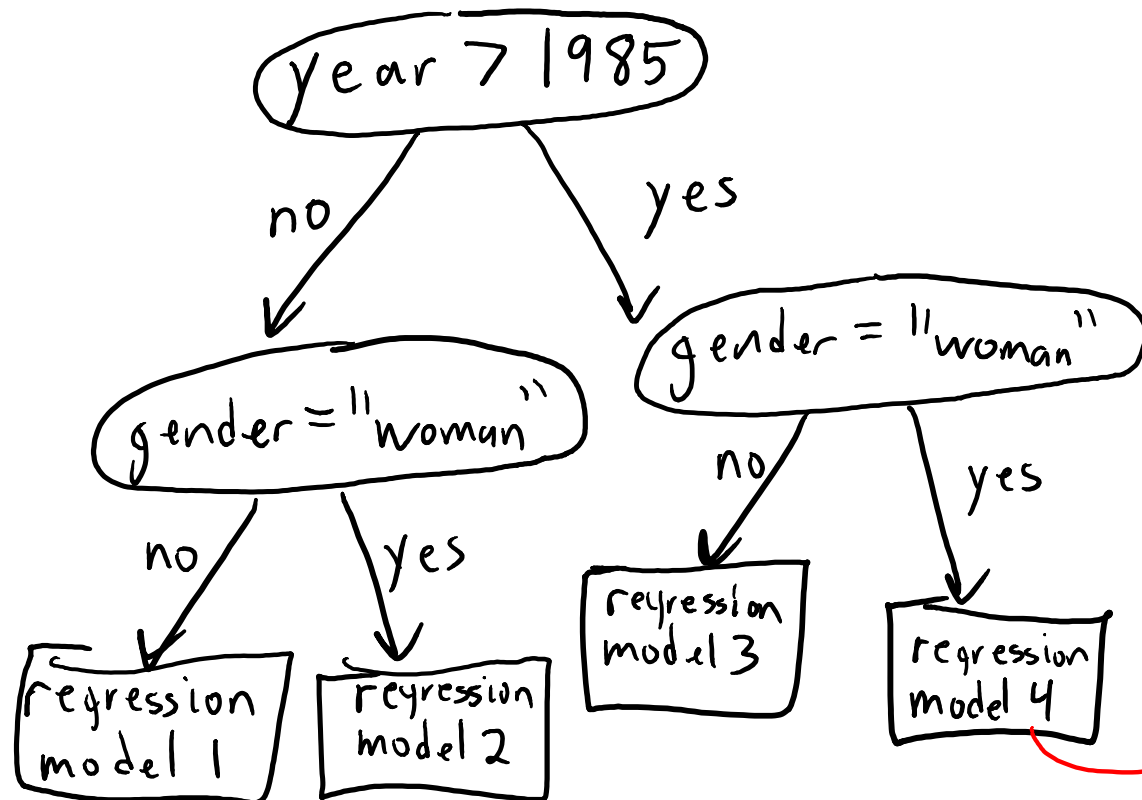SHOT PUT PROGRESSION MEN (7.26 kg) AND WOMEN (4 kg) (mean of top ten)

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  – Regression tree: tree with mean value or linear regression at leaves.

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
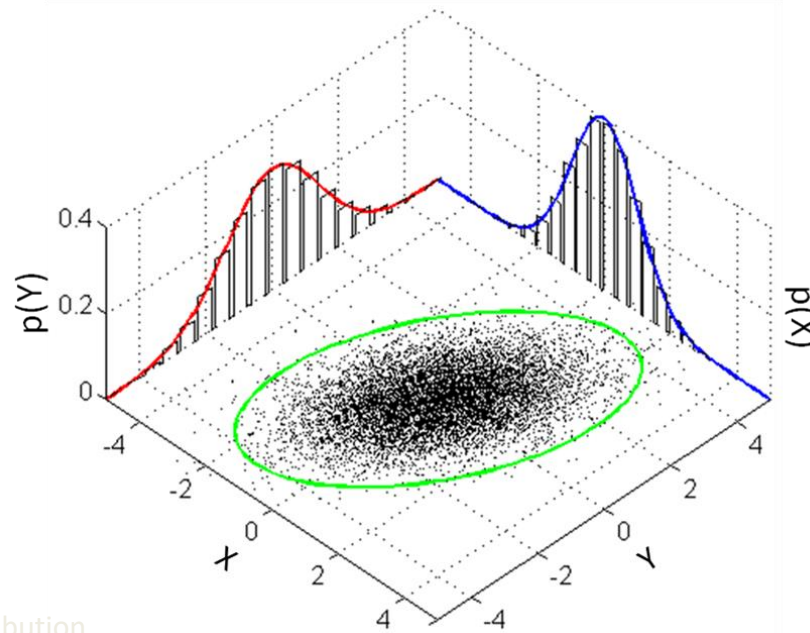  - Regression tree: tree with mean value or linear regression at leaves.
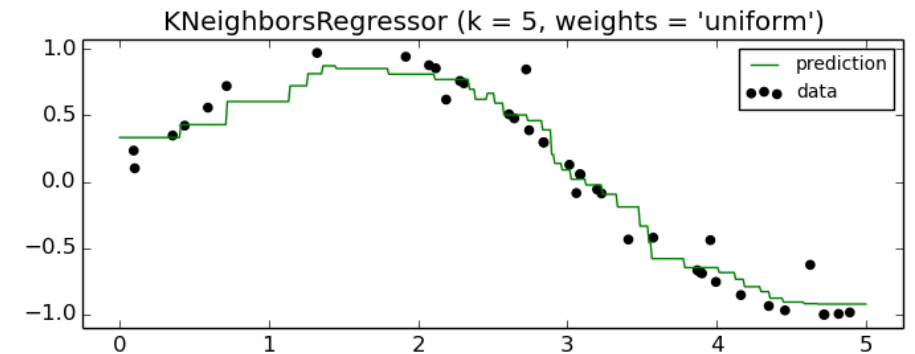  - Generative models: fit $p(x_i \mid y_i)$ and $p(y_i)$ with Gaussian or other model.

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  - Regression tree: tree with mean value or linear regression at leaves.
  - Generative models: fit $p(x_i \mid y_i)$ and $p(y_i)$ with Gaussian or other model.
  - Non-parametric models:
    - Mean $y_i$ among k-nearest neighbours.



KNeighborsRegressor (k = 5, weights = 'uniform')

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  - Regression tree: tree with mean value or linear regression at leaves.
  - Generative models: fit $p(x_i \mid y_i)$ and $p(y_i)$ with Gaussian or other model.
  - Non-parametric models:
    - Mean $y_i$ among k-nearest neighbours.
    - Could be weighted by distance.
      - Close points 'j' get more "weight" $w_{ij}$.
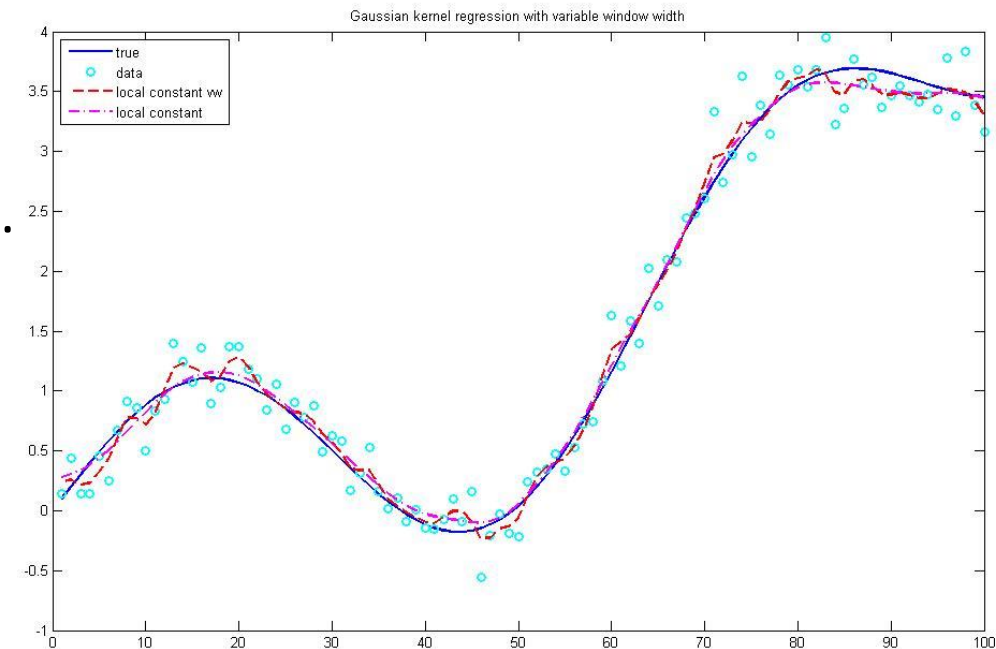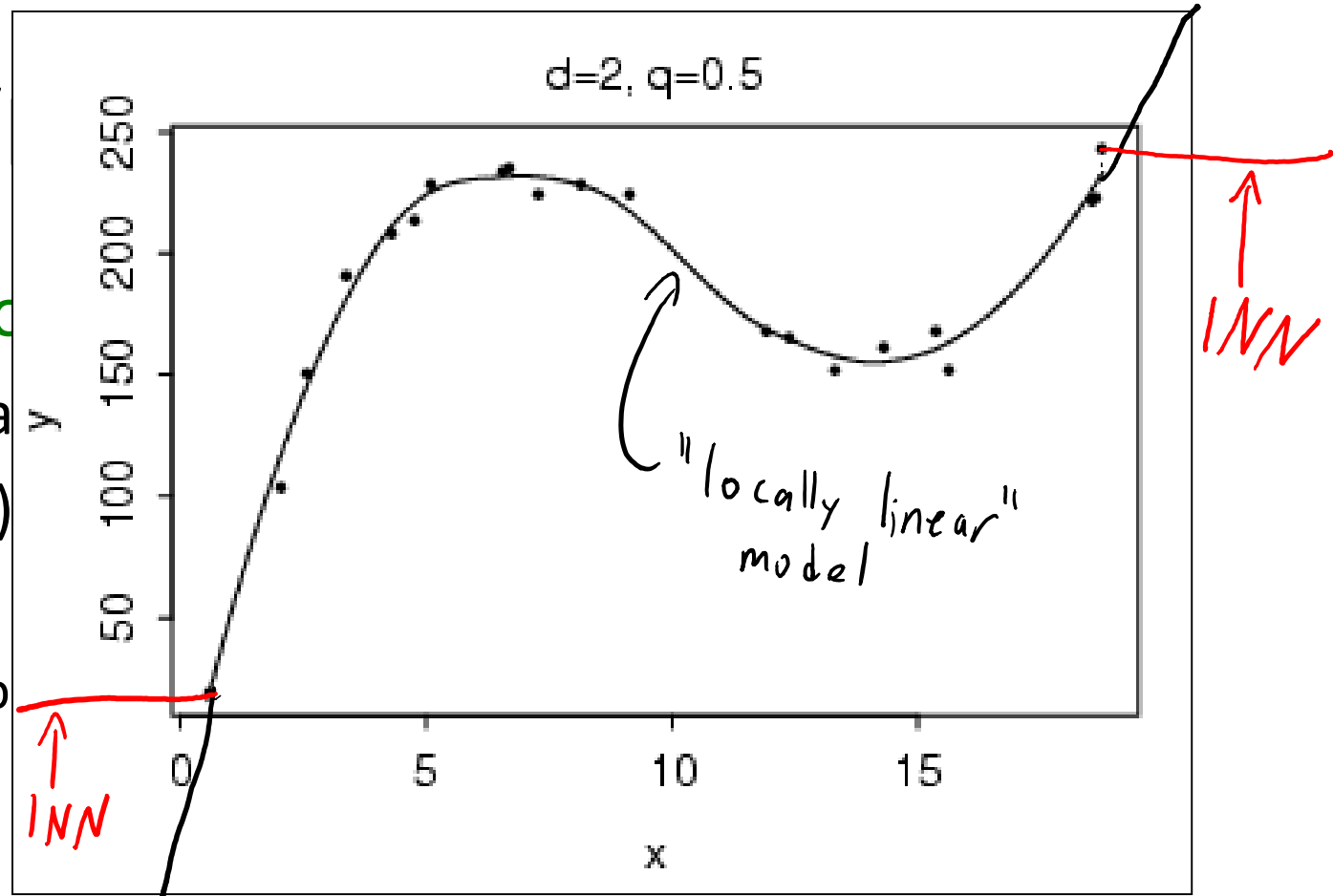
# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  - Regression tree: tree with mean value or linear regression at leaves.
  - Generative models: fit $p(x_i \mid y_i)$ and $p(y_i)$ with Gaussian or other model.
  - Non-parametric models:
    - Mean $y_i$ among k-nearest neighbours.
    - Could be weighted by distance.
    - 'Nadaraya-Waston': weight *all* $y_i$ by distance to $x_i$.

$$y_i = \frac{\sum_{j=1}^{n} w_{ij} y_j}{\sum_{j=1}^{n} w_{ij}}$$

Gaussian kernel regression with variable window width

# Adapting Counting/

- We can adapt our classificatio
  - Regression tree: tree with mea
  - Generative models: fit $p(x_i \mid y_i)$
  - Non-parametric models:
    - Mean $y_i$ among k-nearest neighb
    - Could be weighted by distance.
    - 'Nadaraya-Waston': weight *all* $y_i$
    - 'Locally linear regression': for each $x_i$ a fit linear model weighted by distance.
      (Better than KNN and NW at boundaries.)



d=2, q=0.5

"locally linear" model

1NN

1NN

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  - Regression tree: tree with mean value or linear regression at leaves.
  - Generative models: fit $p(x_i \mid y_i)$ and $p(y_i)$ with Gaussian or other model.
  - Non-parametric models:
    - Mean $y_i$ among k-nearest neighbours.
    - Could be weighted by distance.
    - 'Nadaraya-Waston': weight *all* $y_i$ by distance to $x_i$.
    - 'Locally linear regression': for each $x_i$, fit linear model weighted by distance.
      (Better than KNN and NW at boundaries.)
  - Ensemble methods:
    - Can improve performance by averaging across regression models.

# Regression Forests for Fluid Simulation

- https://www.youtube.com/watch?v=kGB7Wd9CudA

# Linear Least Squares for Quadratic Models

- Can we use linear least squares to fit a quadratic model?

$$y_i = w_0 + w_1 x_i + w_2 x_i^2$$

- You can do this by changing the features (change of basis):

$$X = \begin{bmatrix} 0.2 \\ -0.5 \\ 1 \\ 4 \end{bmatrix} \qquad Z = \begin{bmatrix} 1 & 0.2 & (0.2)^2 \\ 1 & -0.5 & (-0.5)^2 \\ 1 & 1 & (1)^2 \\ 1 & 4 & (4)^2 \end{bmatrix}$$
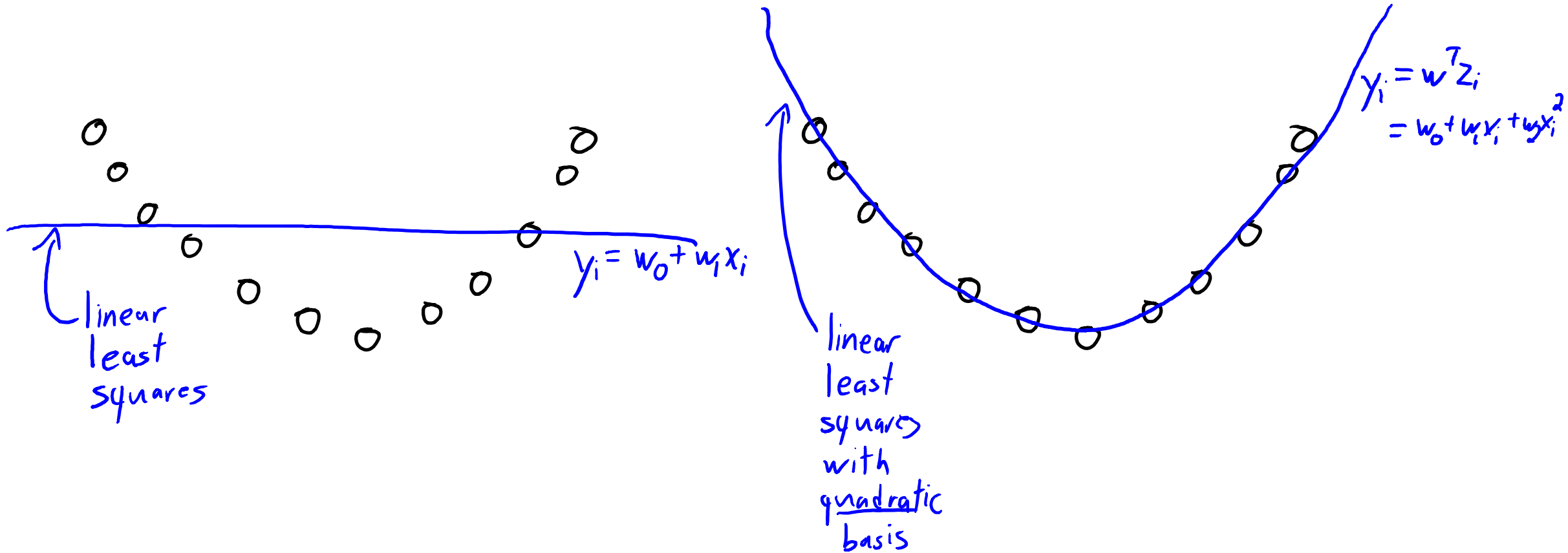
$$\underset{y-\text{inf}}{\phantom{1}} \quad \underset{x}{\phantom{0.2}} \quad \underset{x^2}{\phantom{(0.2)^2}}$$

$$y_i = w^\top z_i$$
$$= w_0 z_{i0} + w_1 z_{i1} + w_2 z_{i2}$$
$$\approx w_0 + w_1 x_i + w_2 x_i^2$$

- It's a linear function of w, but a quadratic function of $x_i$.

- Fitting with least squares: $w = (Z^\top Z) \backslash (Z^\top y)$

To predict on new data $\hat{X}$, form $\hat{Z}$ from $\hat{X}$ and take $y = \hat{Z} w$

# Linear Least Squares for Quadratic Models


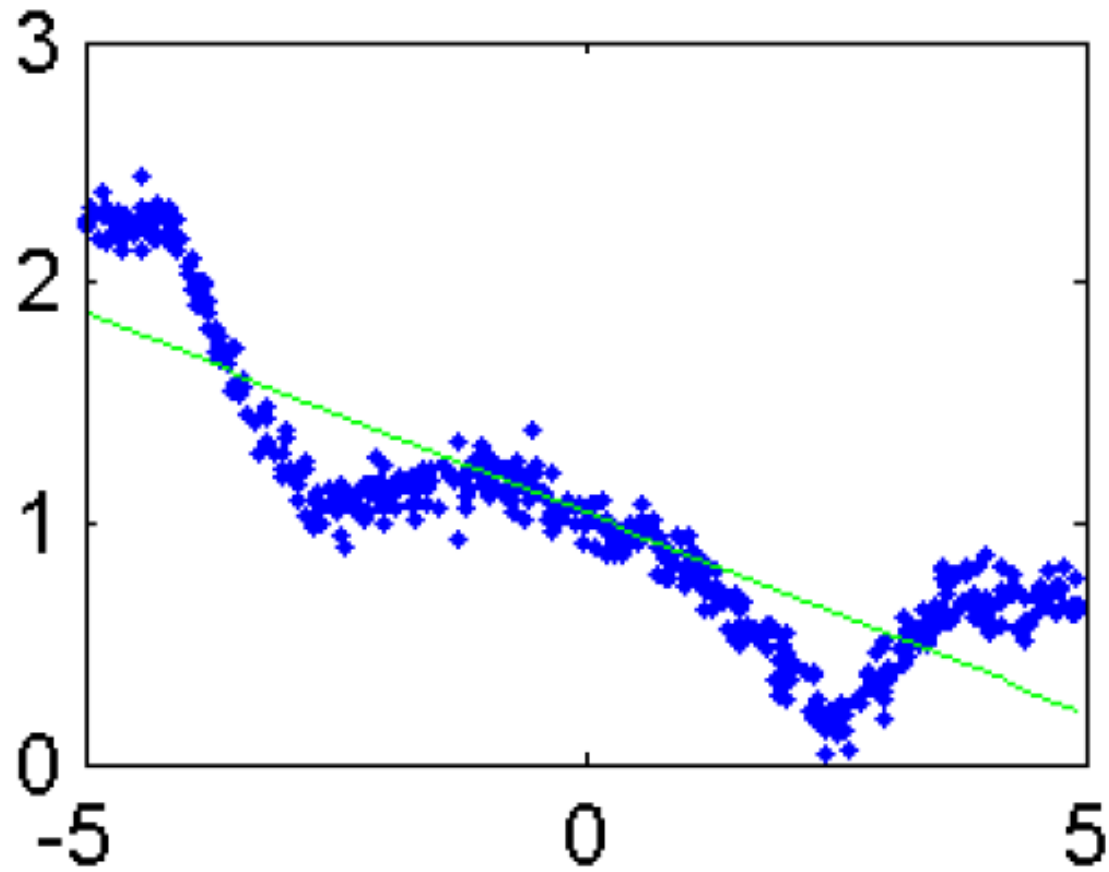
$y_i = w_0 + w_1 x_i$

linear least squares

$y_i = w^T z_i$
$= w_0 + w_1 x_i + w_2 x_i^2$

linear least squares with quadratic basis

# General Polynomial Basis

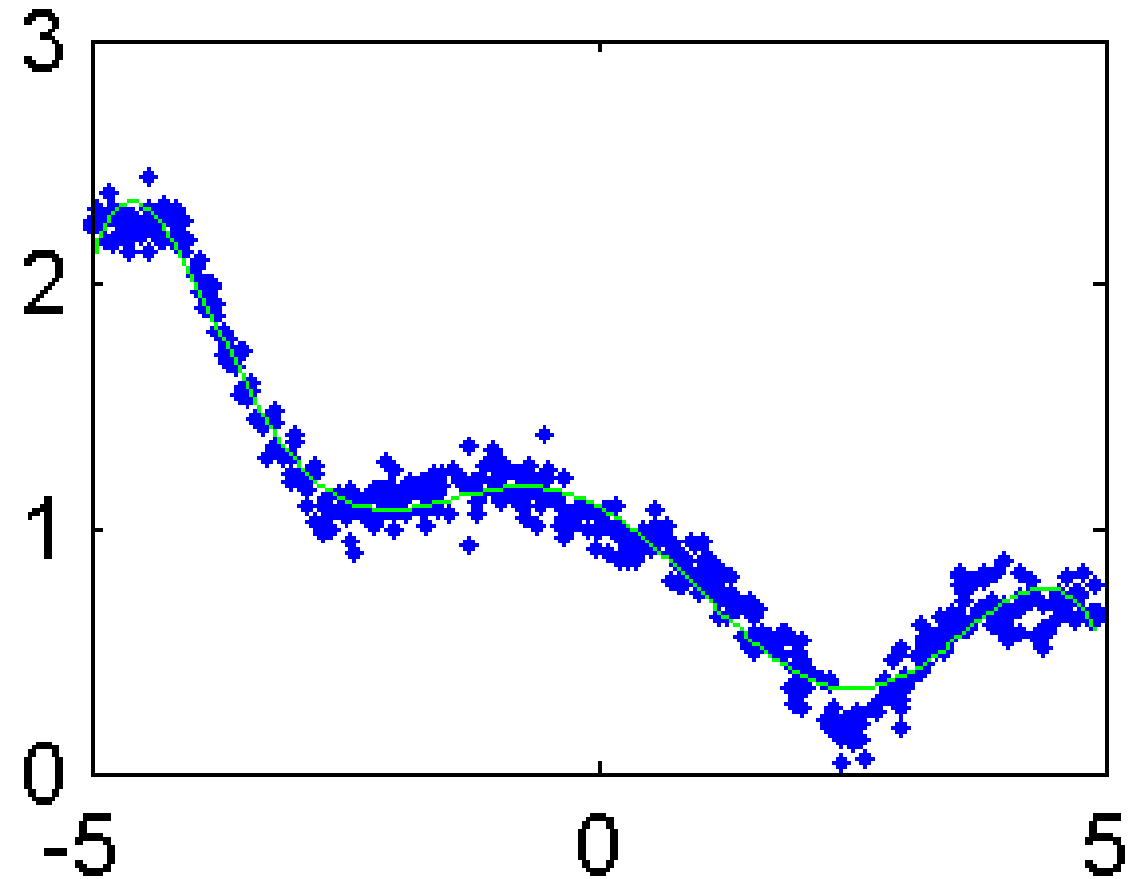- We can have a polynomial of degree 'p' by using a basis:

$$Z = \begin{bmatrix} 1 & x_1 & (x_1)^2 & \text{----} & (x_1)^p \\ 1 & x_2 & (x_2)^2 & \text{----} & (x_2)^p \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & (x_n)^2 & \text{----} & (x_n)^p \end{bmatrix}$$

- There are polynomial basis functions that are numerically nicer:
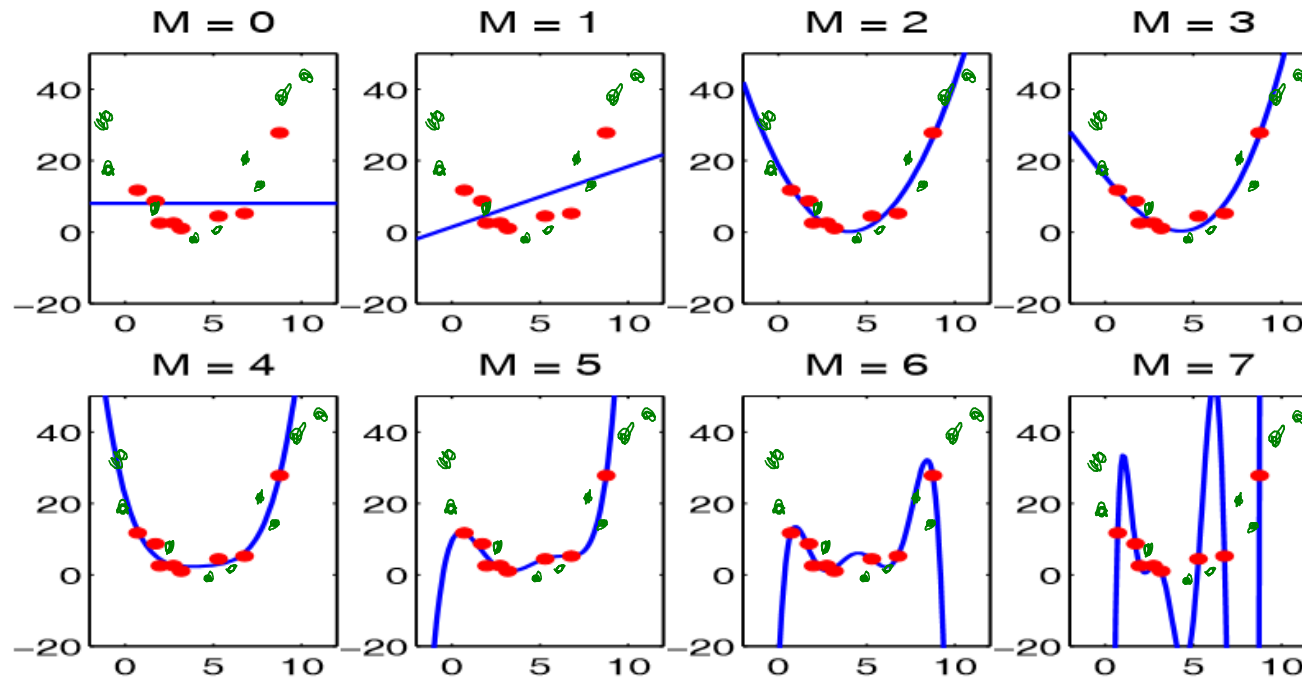  - E.g., Lagrange polynomials.

# General Polynomial Basis



Degree 7

# Degree of Polynomial and Fundamental Trade-Off

- As the polynomial degree increases, the <span style="color:red">training error</span> goes down.



- But training error becomes worse approximation <span style="color:green">test error</span>.
- Usual approach to <span style="color:green">selecting degree</span>: <span style="color:blue">validation</span> or <span style="color:blue">cross-validation</span>.

# Summary

- Y-intercept can be modeled by using a column of 1s.
- Linear least squares solution is given by normal equations:
    - Solve $(X^TX)w = X^Ty$.
- Tree/generative/non-parametric/ensemble methods for regression.
- Change of basis allows linear models to model non-linear data:


- Next time:
    - Bases that can model any continuous function.

# Bonus Slide: Householder(-ish) Notation

- **Househoulder notation:** set of (fairly-logical) conventions for math.

Use greek letters for scalars: $\alpha = 1$, $\beta = 3.5$, $\gamma = \pi$

Use first/last lowercase letters for vectors: $w = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$, $x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $y = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$, $a = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$, $b = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

$\longrightarrow$ Assumed to be column-vectors.

Use first/last uppercase letters for matrices: $X$, $Y$, $W$, $A$, $B$

Indices use $i, j, k$.

Sizes use $m, n, d, p,$ and $k$ $\longleftarrow$ hopefully meaning of 'k' is obvious from context

Sets use $S, T, U, V$

Functions use $f, g,$ and $h$.

When I write $x_i$ I mean "grab row 'i' of X and make a column-vector with its values."

# Bonus Slide: Householder(-ish) Notation

- **Househoulder notation:** set of (fairly-logical) conventions for math:

Our ultimate least squares notation:

$$f(w) = \frac{1}{2} \| Xw - y \|^2$$

But if we agree on notation we can quickly understand:

$$g(x) = \frac{1}{2} \| Ax - b \|^2$$

If we use random notation we get things like:

$$H(\beta) = \frac{1}{2} \| R\beta - P_n \|^2$$

Is this the same model?