

CPSC 340: Machine Learning and Data Mining

Non-Linear Regression

Admin

- Assignment 1 grades are out.
- Assignment 2 is due Sunday.
- Assignment 3 will be out by early next week.
 - Will be due before the break
- Midterm is after the break (March 1 in class)
- Tutorial next week: practice problems for hw3

Last Time: Linear Regression

- We discussed **linear models**:

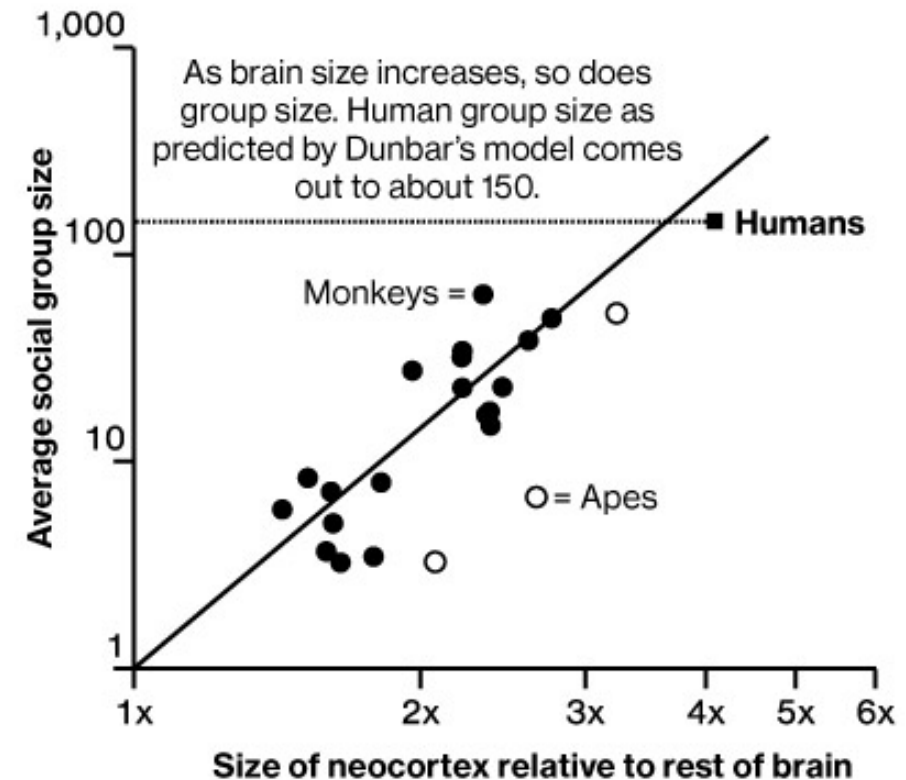
$$y_i = w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} \\ = \sum_{j=1}^d w_j x_{ij} = w^T x_i$$

- “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 \leftarrow $w^T x_i$ \rightarrow True value y_i

The Social Cortex



DATA: THE SOCIAL BRAIN HYPOTHESIS, DUNBAR 1998

To predict on test case \hat{x}_i
use $\hat{y}_i = w^T \hat{x}_i$

Last Time: Supervised Learning Notation

- We're treating 'w', 'y', and each x_i as **column-vectors**:

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}_{d \times 1}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

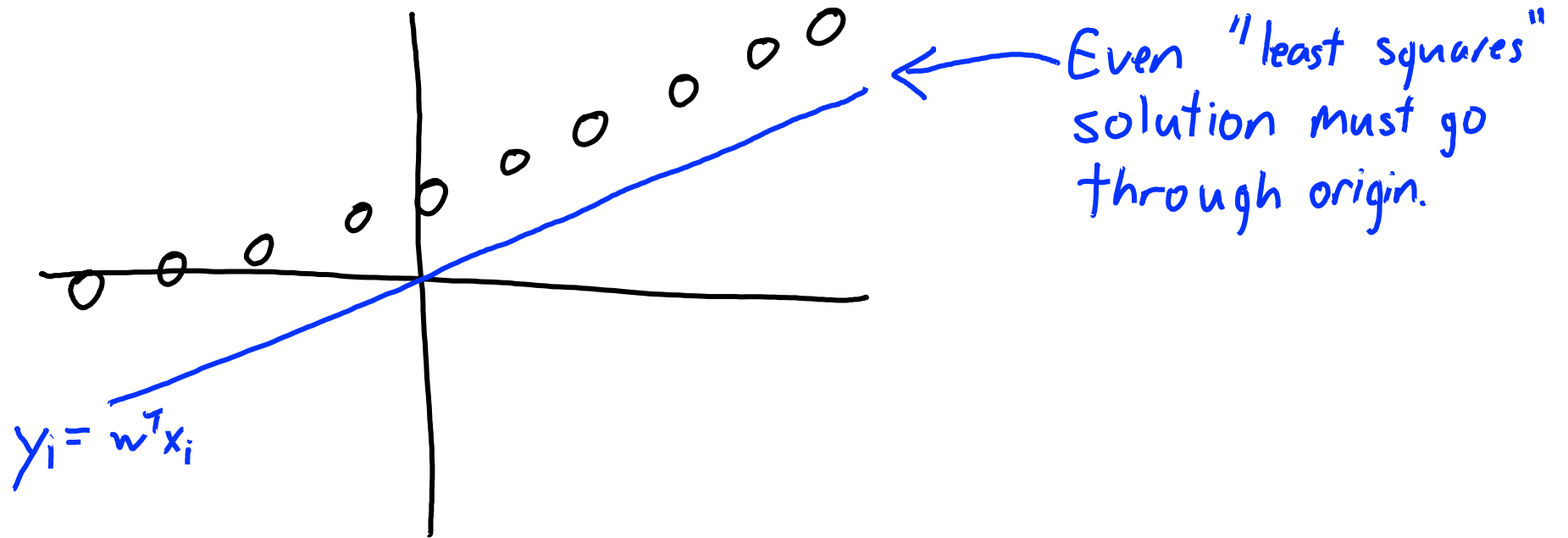
$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}_{d \times 1}$$

- So feature matrix 'X' actually has x_i **transposed as rows**:

$$X = \begin{bmatrix} \text{---} x_1^T \text{---} \\ \text{---} x_2^T \text{---} \\ \vdots \\ \text{---} x_n^T \text{---} \end{bmatrix}$$

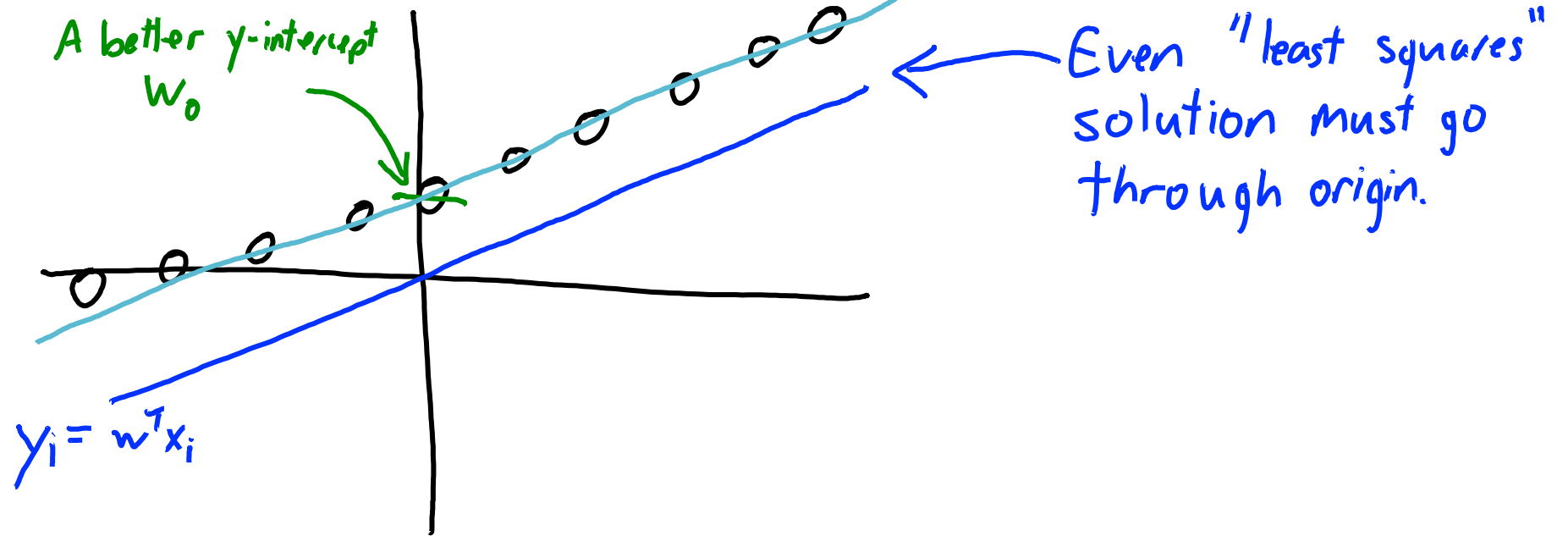
Why don't we have a y-intercept?

- 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 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”**.

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

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

(Note: A red bracket and 'X' mark are drawn under the last two columns of Z, indicating they are the original features X.)

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

$$\begin{aligned} y_i &= w_0 z_{i0} + w_1 z_{i1} + w_2 z_{i2} \\ &\quad \quad \quad \downarrow \quad \quad \downarrow \quad \quad \downarrow \\ &\quad \quad \quad \text{"1"} \quad \quad x_{i1} \quad \quad x_{i2} \\ &= w_0 + w_1 x_{i1} + w_2 x_{i2} \end{aligned}$$

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

Linear Least Squares

Prediction:

$$y = X * w$$

Why?

$$y_i = w^T x_i \quad \text{so}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_n \end{bmatrix} = \begin{bmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{bmatrix} = \underbrace{\begin{bmatrix} \text{---} x_1^T \text{---} \\ \text{---} x_2^T \text{---} \\ \vdots \\ \text{---} x_n^T \text{---} \end{bmatrix}}_X \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = Xw$$

Linear Least Squares

- We can rewrite the objective function as a norm

Want ' w ' that minimizes

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

Define "residual" r_i as

Signed error on example ' i ':

$$r_i = w^T x_i - y_i$$

$$r = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{bmatrix} = \begin{bmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_n - y_n \end{bmatrix} = \underbrace{\begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_n \end{bmatrix}}_{Xw} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \boxed{Xw - y}$$

Linear Least Squares

Want ' w ' that minimizes

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

$$\boxed{\frac{1}{2} \|Xw - y\|^2}$$

Let's expand
then compute
gradient.

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

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

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

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

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

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

A good way to check
your steps: make sure
that dimensions all make sense.

Linear Least Squares

See notes on linear and quadratic derivatives for details.

Training: $w = \text{solve}(X.T @ X, X.T @ y)$

Why?

Want 'w' that minimizes

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

a quadratic function (in matrix notation)

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

What are the gradients of these terms?

So at a minimizer where $\nabla f(w) = 0$
we have: $X^T X w = X^T y$ "normal equations"

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

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

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

This is like saying $\frac{d}{dw} [\alpha w] = \alpha$

Some matrix Some vector
This is a linear system $Aw = b$

for some matrix 'A' and vector 'b'

The Punch Line

$$\min_w \frac{1}{2} \|Xw - y\|_2^2$$



$$w = \text{solve}(X.T @ X, \quad X.T @ y)$$

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

Incorrect Solutions to Least Squares Problem

The least squares objective is $f(w) = \frac{1}{2} \|Xw - y\|^2$

The minimizers of this objective are solutions to the linear system:

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

The following are not the solutions to the least squares problem:

$$w = (X^T X)^{-1} (X^T y) \quad (\text{only true if } \underline{X^T X \text{ is invertible}})$$

$$w X^T X = X^T y \quad (\text{matrix multiplication is } \underline{\text{not}} \text{ commutative, dimensions don't even match})$$

$$w = \frac{X^T y}{X^T X} \quad (\text{you } \underline{\text{cannot divide by a matrix}})$$

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 be so big we **can't store $X^T X$** .
 - 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 $X^T X$ is $d \times d$.

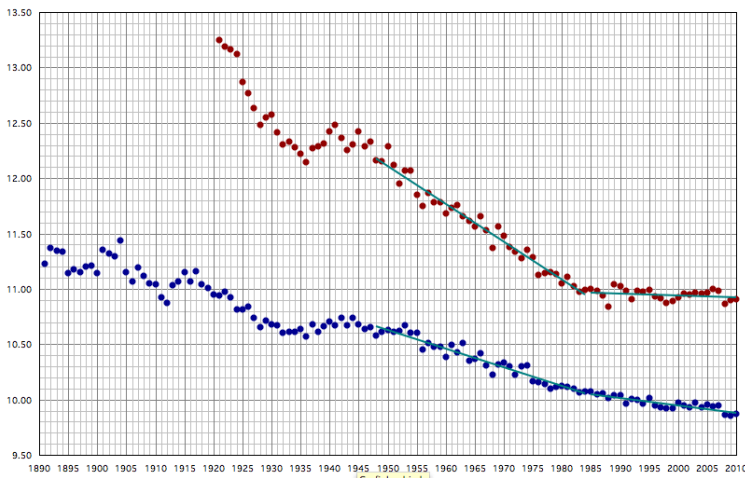
Costs $O(nd^2)$ to calculate:

- Each of the $O(d^2)$ elements is an inner product between length ' n ' vectors.

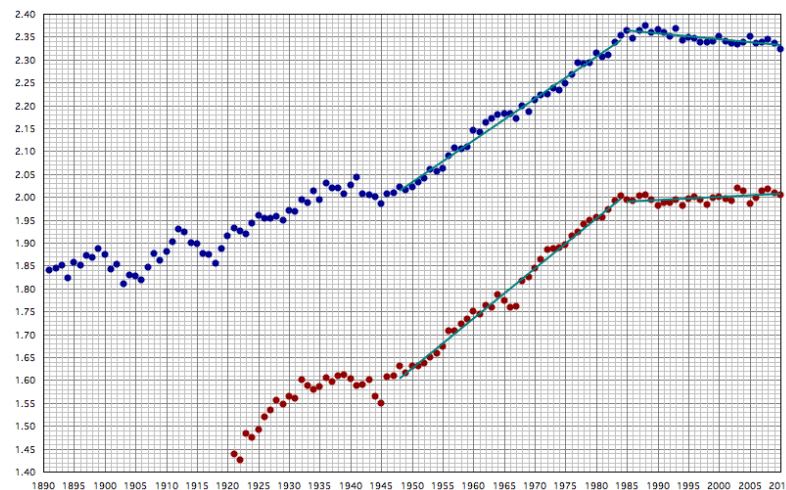
Example: Non-Linear Progressions in Athletics

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

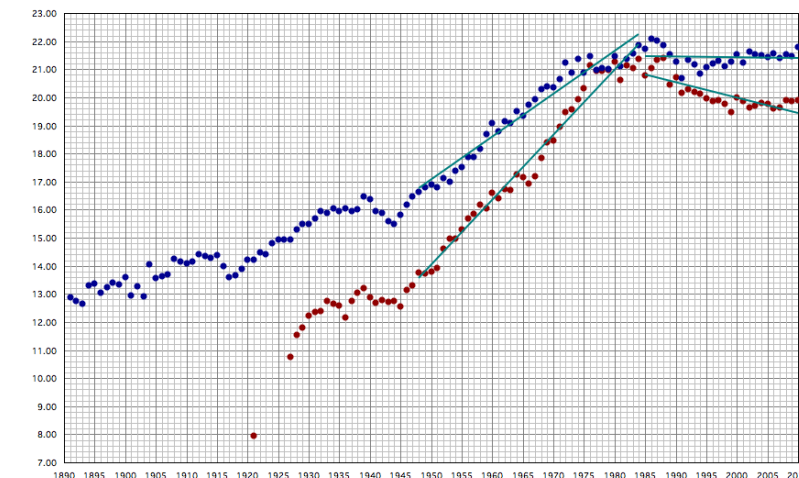
100m PROGRESSION MEN AND WOMEN (mean of top ten)



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)

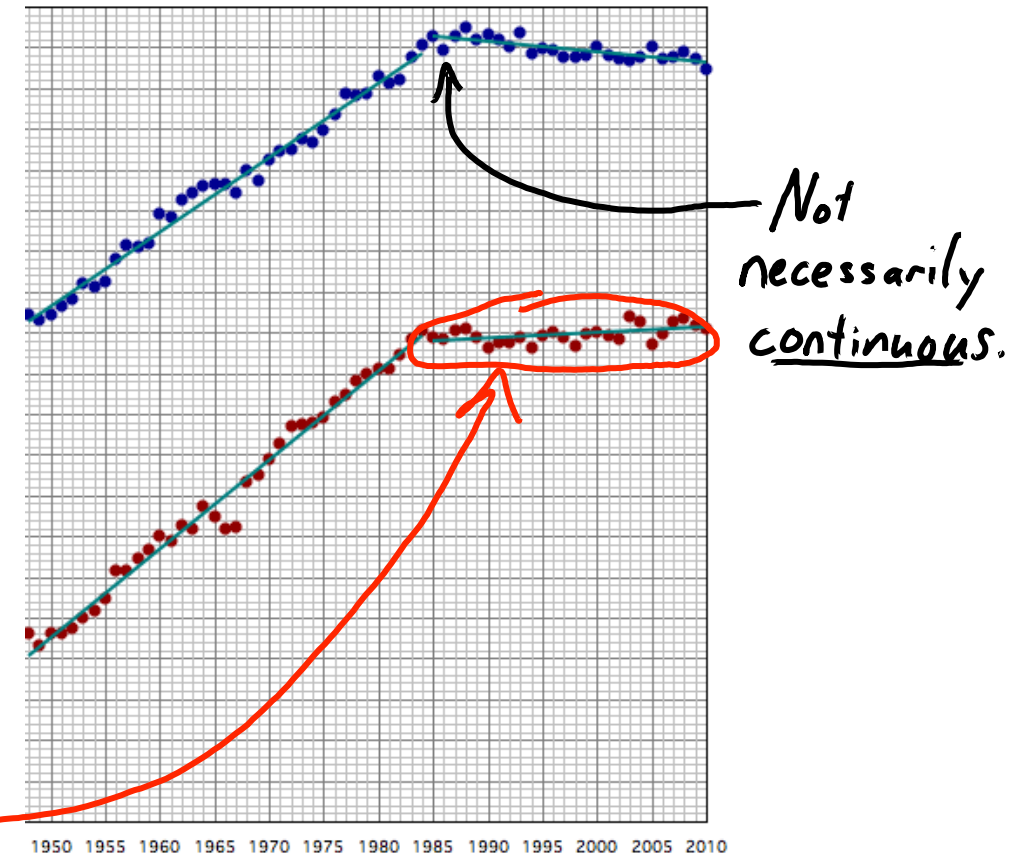
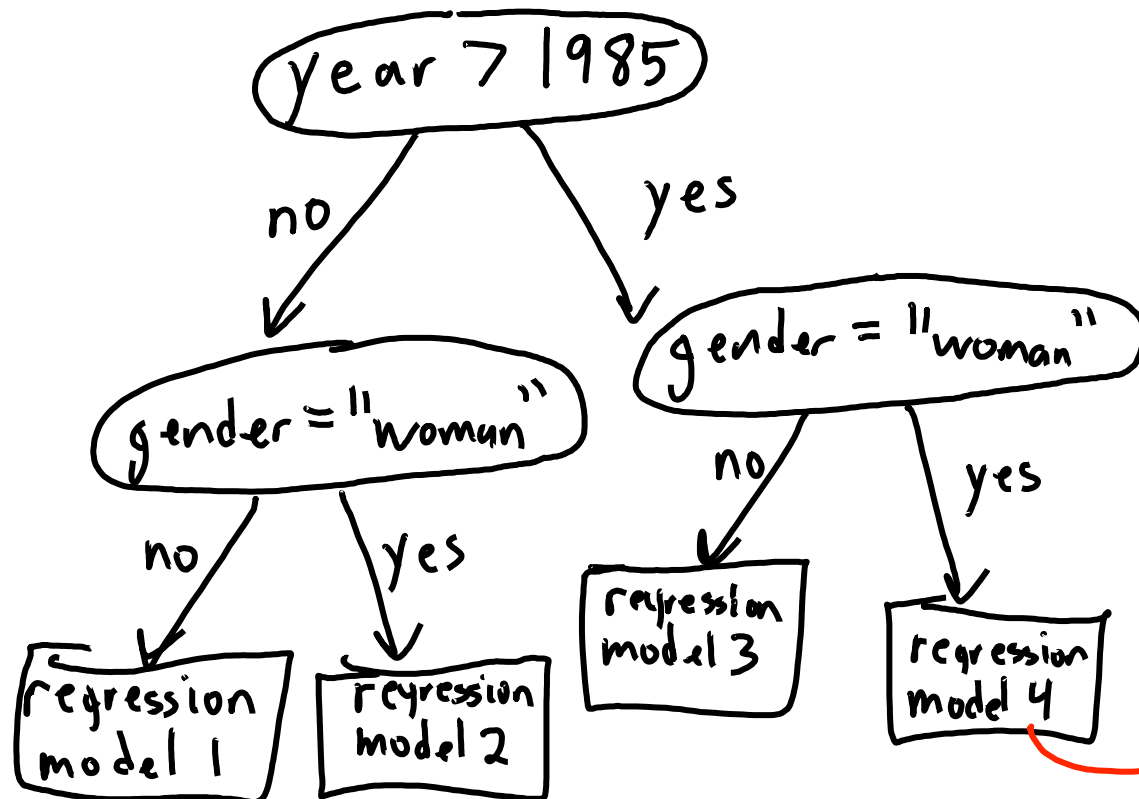


Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:

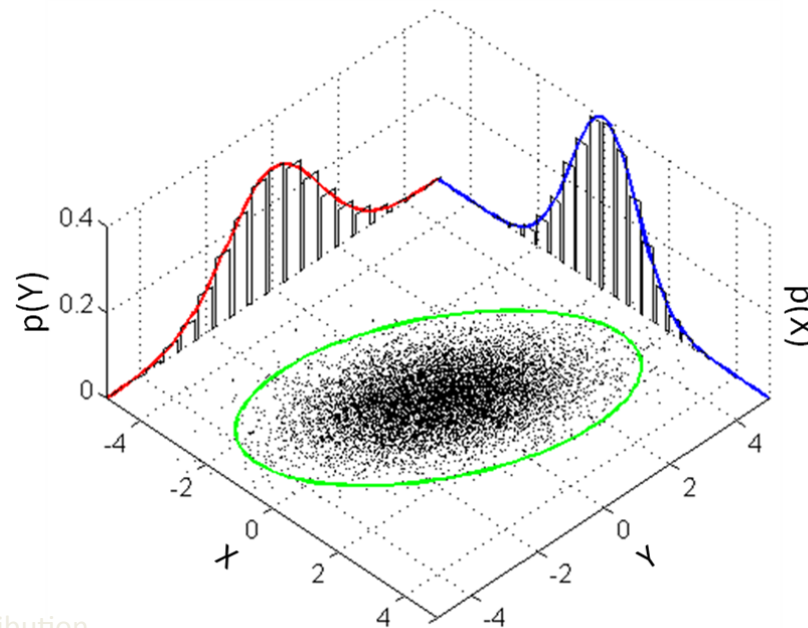
Adapting Counting/Distance-Based Methods

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



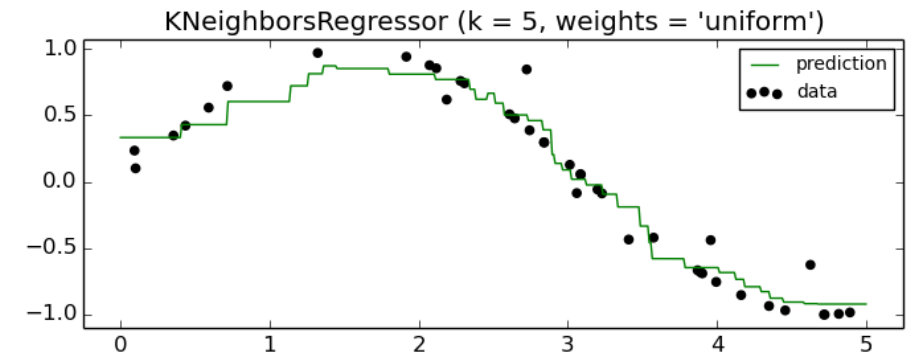
Adapting Counting/Distance-Based Methods

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



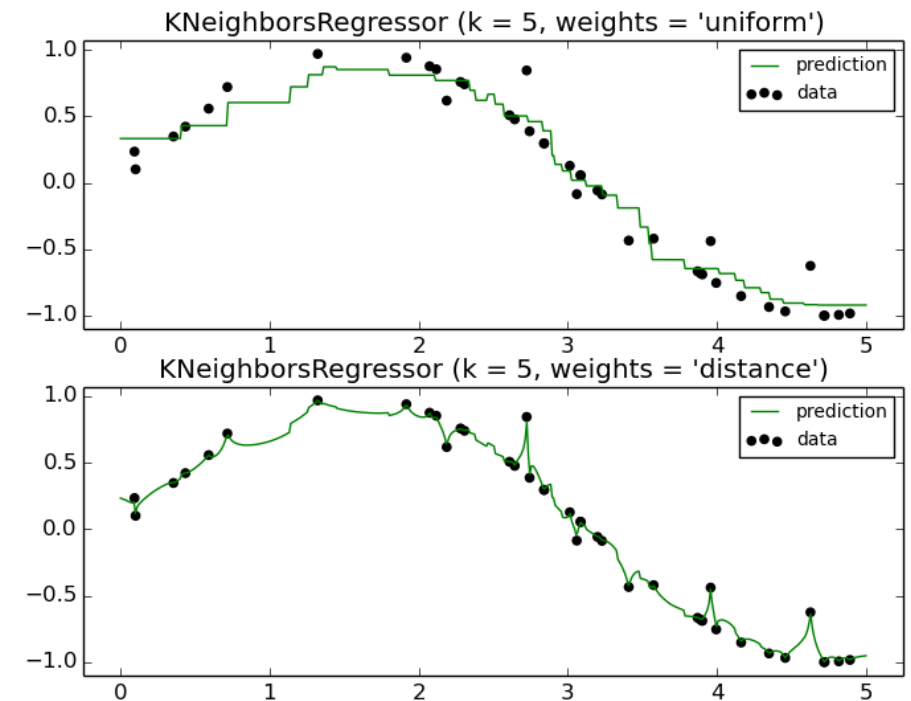
Adapting Counting/Distance-Based Methods

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



Adapting Counting/Distance-Based Methods

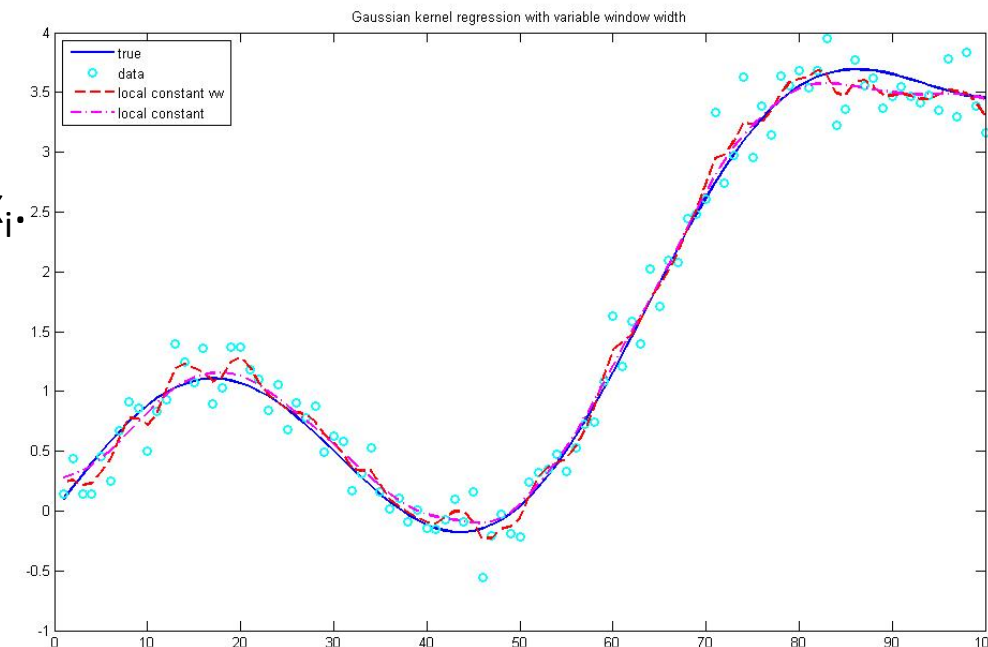
- We can adapt our classification methods to perform regression:
 - Regression tree: tree with mean value or linear regression at leaves.
 - Generative models: fit $p(x_i | y_i)$ and $p(y_i)$ with Gaussian or other model.
 - Non-parametric models:
 - Mean y_i among k -nearest neighbours.
 - Could be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

$$\hat{y}_i = \frac{\sum_{j=1}^n w_{ij} y_j}{\sum_{j=1}^n w_{ij}}$$



Adapting Counting/

- We can **adapt our classification**

- Regression tree: tree with mean

- Generative models: fit $p(x_i | y_i)$

- Non-parametric models:

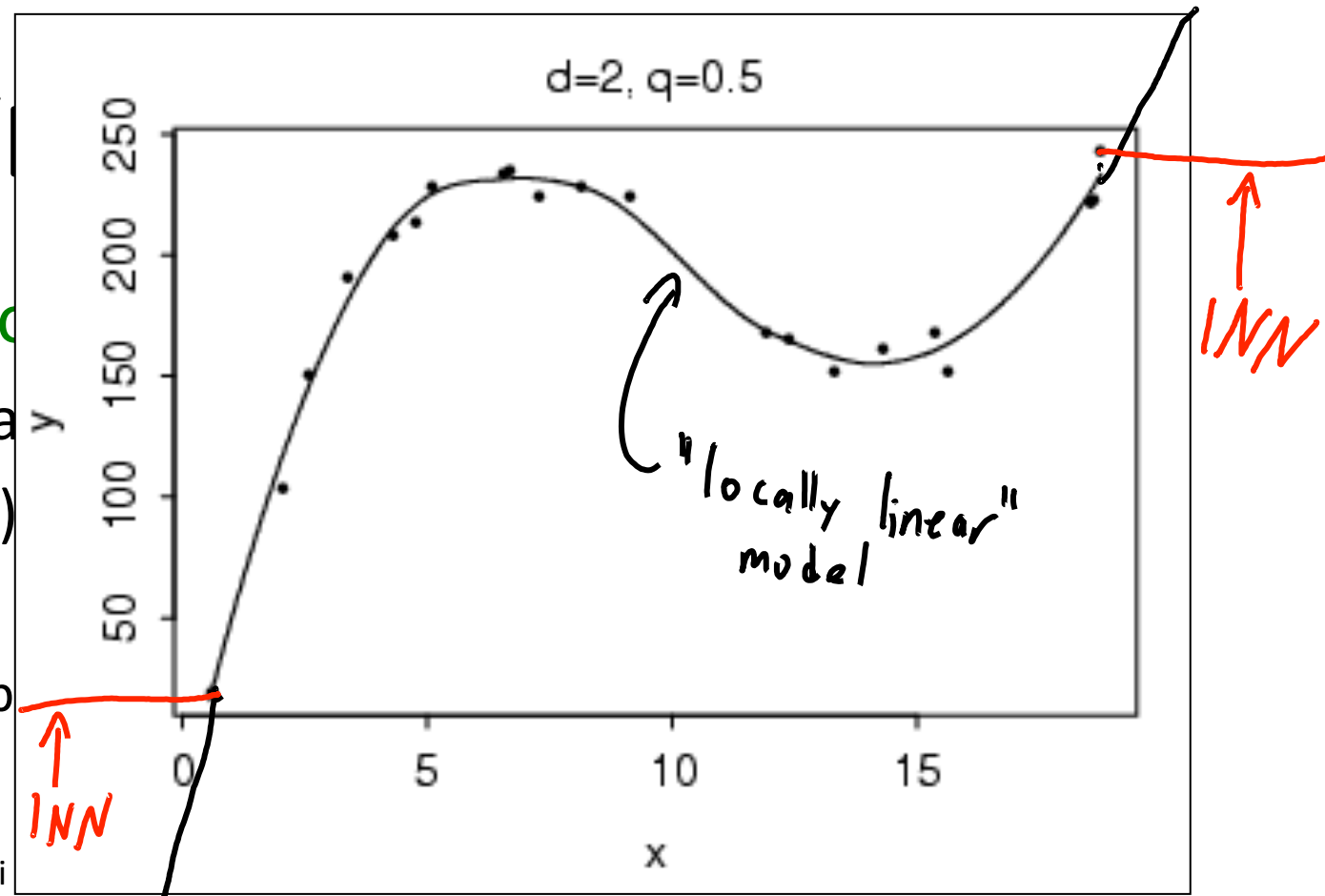
- Mean y_i among k-nearest neighbors

- 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.)



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 $p(x_i | y_i)$ 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.

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**):

$$X = \begin{bmatrix} 0.2 \\ -0.5 \\ 1 \\ 4 \end{bmatrix}$$

$$Z = \begin{bmatrix} 1 & 0.2 & (0.2)^2 \\ 1 & -0.5 & (-0.5)^2 \\ 1 & 1 & (1)^2 \\ 1 & 4 & (4)^2 \end{bmatrix}$$

y -int
 x
 x^2

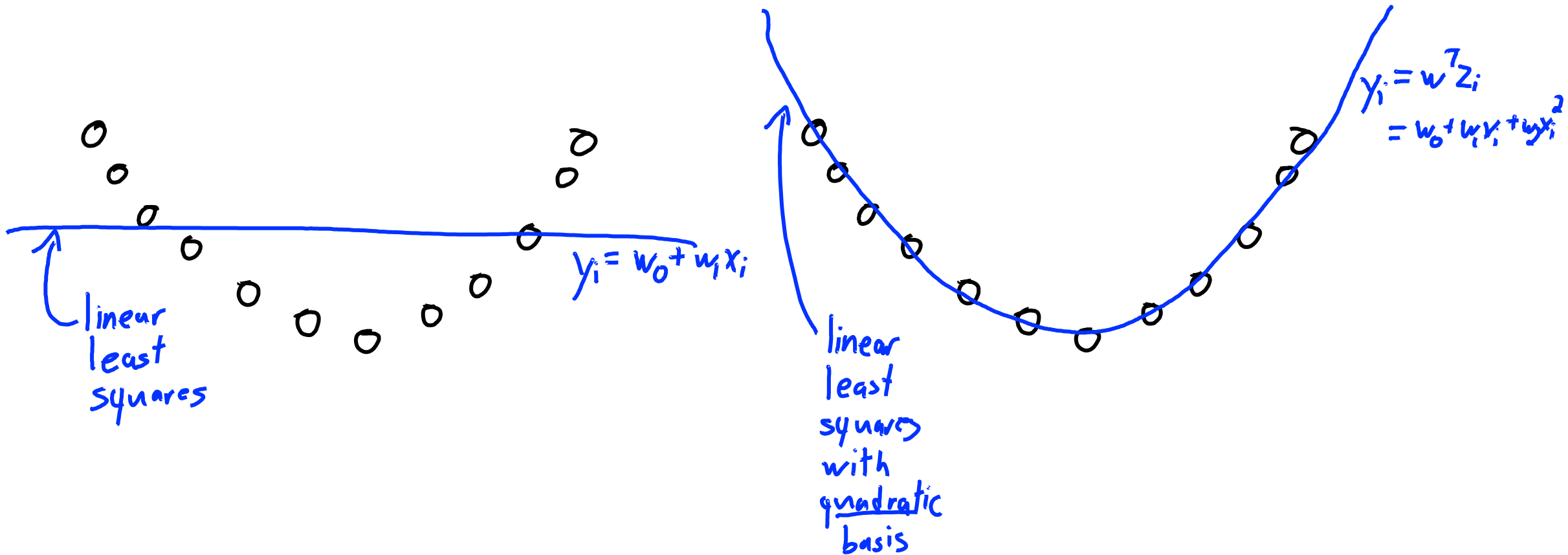
$$\begin{aligned} y_i &= w^T z_i \\ &= w_0 z_{i0} + w_1 z_{i1} + w_2 z_{i2} \\ &= w_0 + w_1 x_i + w_2 x_i^2 \end{aligned}$$

"solve
linear
system"

- Fitting with least squares: $w = (Z^T Z)^{-1} (Z^T y)$
- It's a **linear function of w** , but a **quadratic function of x_i** .

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



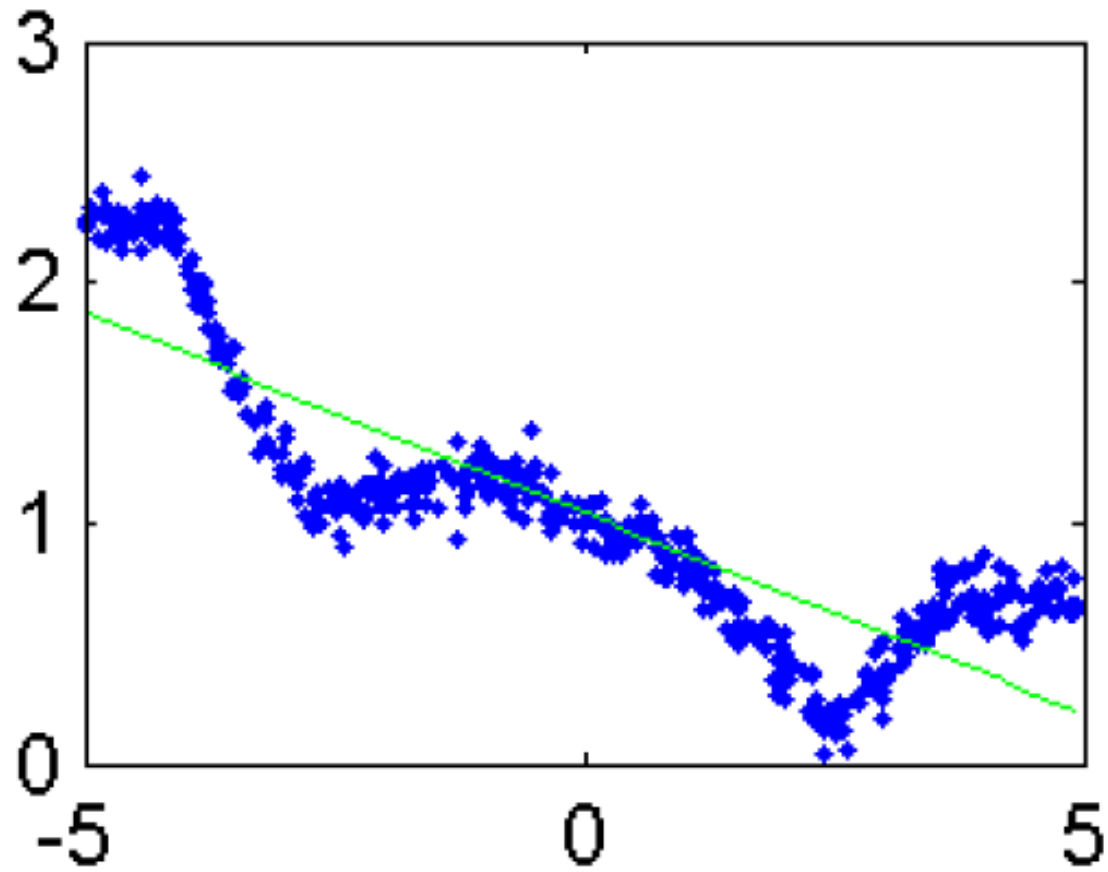
General Polynomial Basis

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

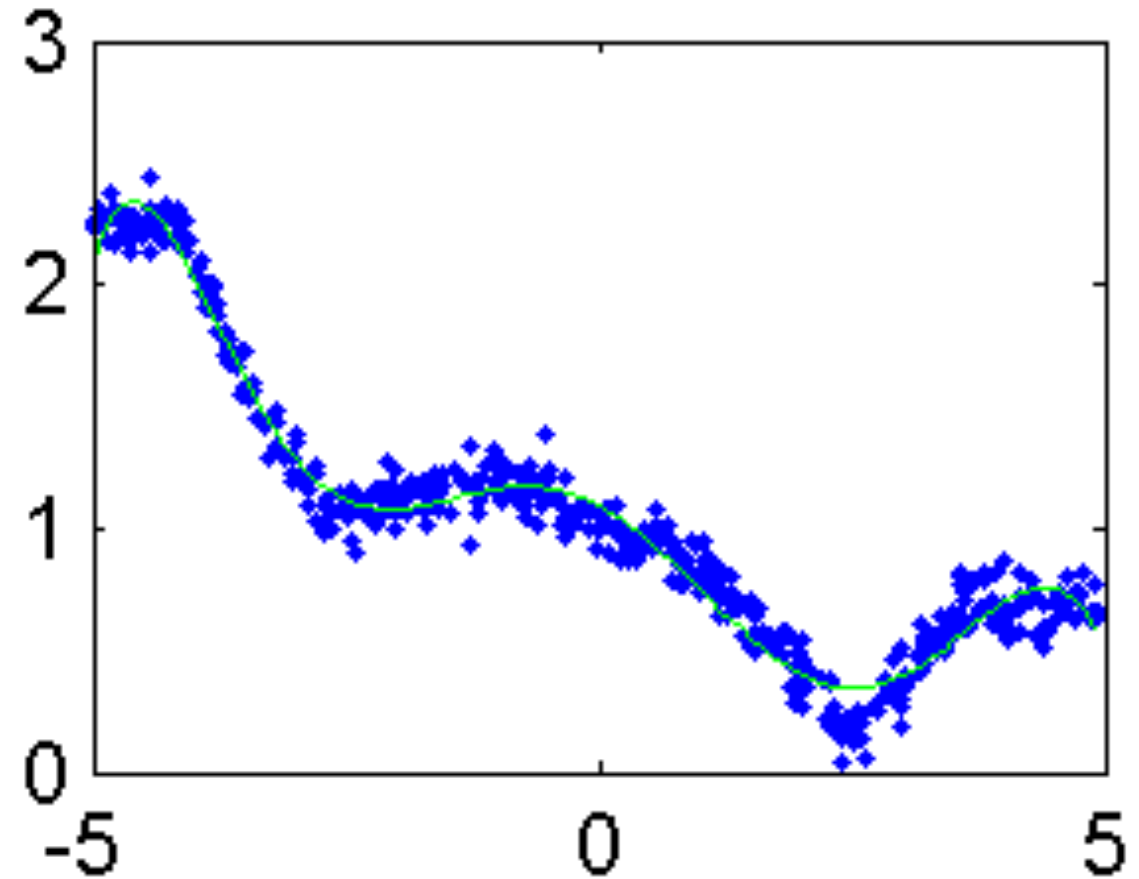
$$Z = \begin{bmatrix} 1 & x_1 & (x_1)^2 & \dots & (x_1)^p \\ 1 & x_2 & (x_2)^2 & \dots & (x_2)^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & (x_n)^2 & \dots & (x_n)^p \end{bmatrix}$$

- There are polynomial basis functions that are numerically nicer:
 - E.g., Lagrange polynomials (see CPSC 303)

General Polynomial Basis

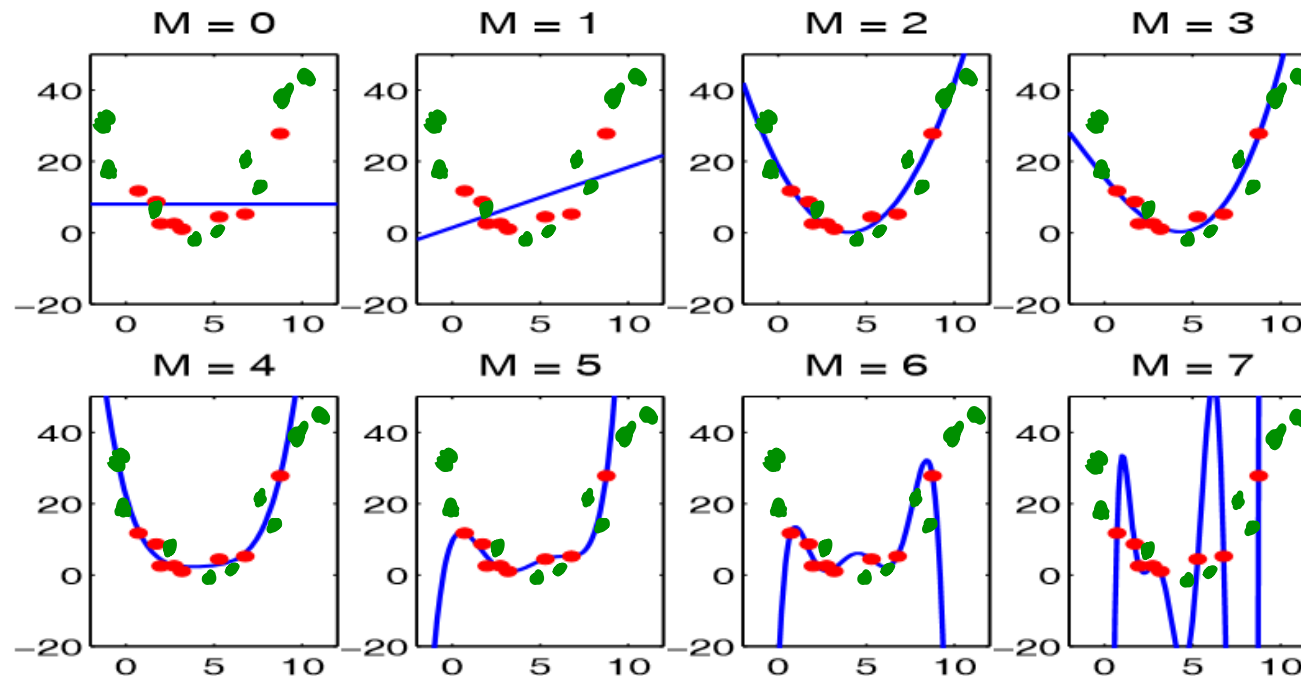


Degree 7



Degree of Polynomial and Fundamental Trade-Off

- As the polynomial degree increases, the **training error** goes down.



- But training error becomes worse approximation **test error**.
- Usual approach to **selecting degree**: **validation** or **cross-validation**.

Summary

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