# CPSC 340:
# Machine Learning and Data Mining

Non-Linear Regression

# Admin

- <span style="color:red">Assignment 1</span> grades are out.
- <span style="color:red">Assignment 2</span> is due Sunday.
  - Extra office hours added on Saturday 12-2pm (see office hours calendar)
- <span style="color:red">Assignment 3</span> will be out by early next week.
  - Will be due before the break
- <span style="color:red">Midterm</span> is after the break (March 1 in class)
- Tutorial next week: practice problems for hw3

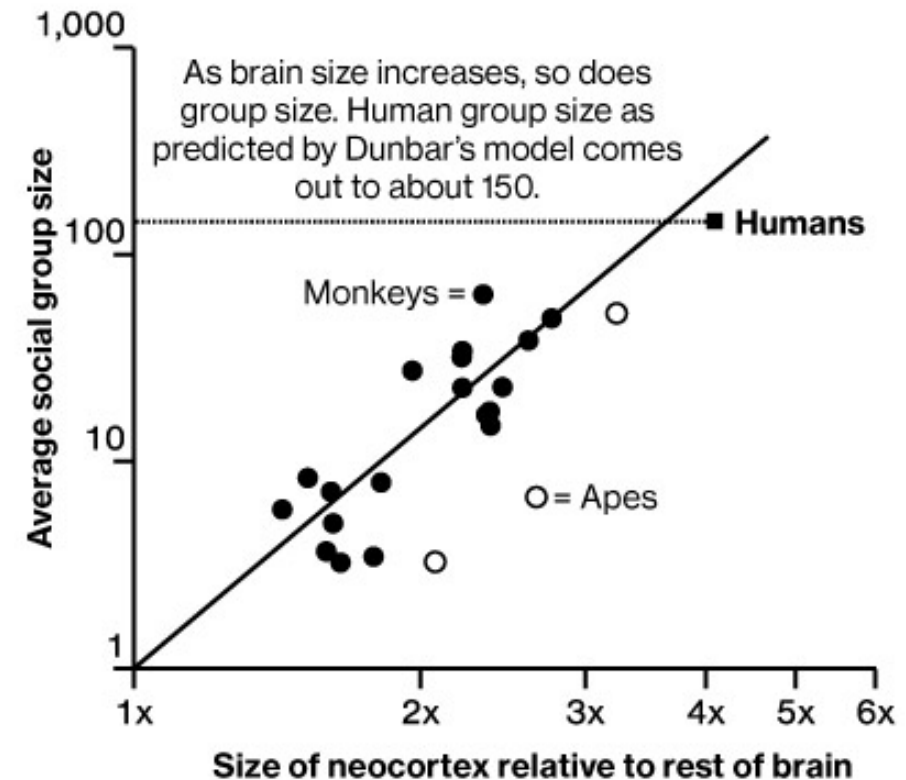# Last Time: Linear Regression

- We discussed linear models:

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

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

- "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

**The Social Cortex**

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

■ Humans

Monkeys = ●

O = Apes

Average social group size

1,000 / 100 / 10 / 1

Size of neocortex relative to rest of brain

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

DATA: THE SOCIAL BRAIN HYPOTHESIS, DUNBAR 1998

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

3

# Last Time: Supervised Learning Notation

- We're treating 'w', 'y', and each $x_i$ as <span style="color:blue">column-vectors</span>:

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_n \end{bmatrix} \quad x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ \vdots \\ x_{id} \end{bmatrix}$$
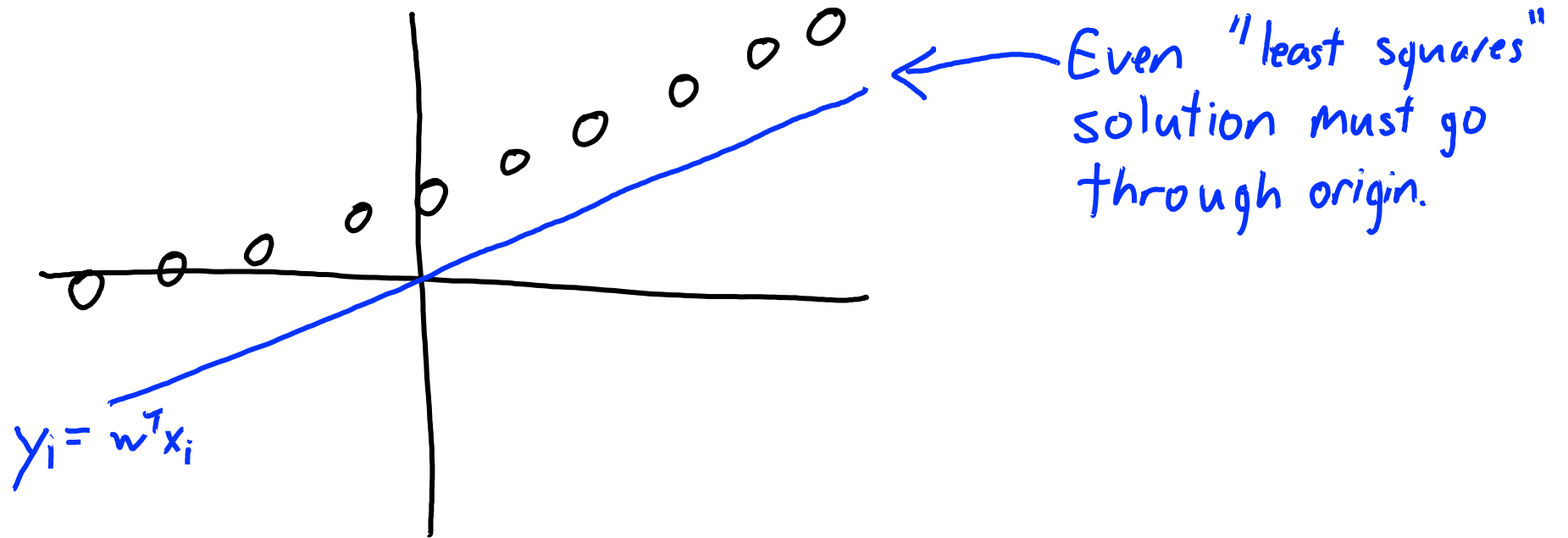
$$d \times 1 \quad\quad n \times 1 \quad\quad d \times 1$$

- So feature matrix 'X' actually has <span style="color:green">$x_i$ transposed as rows</span>:

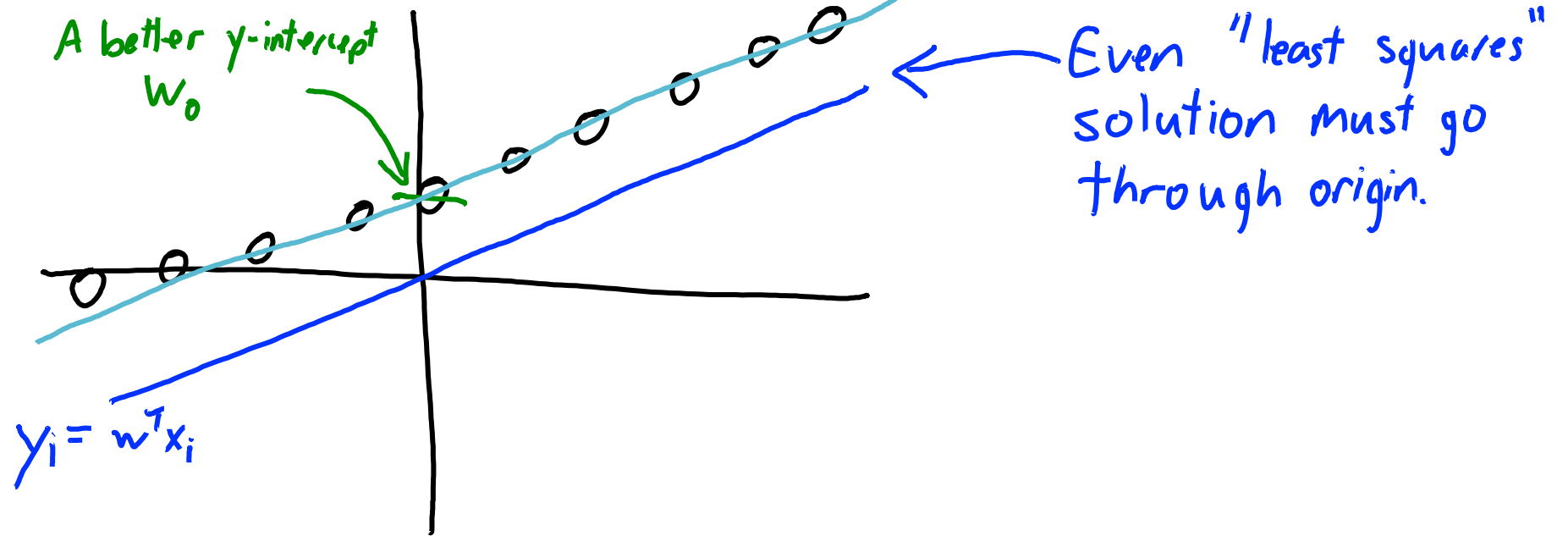$$X = \begin{bmatrix} - x_1^7 - \\ - x_2^7 - \\ \vdots \\ - x_n^7 - \end{bmatrix}$$

# 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 $y_i = 0$.

$y_i = w^T x_i$

Even "least squares" solution must go through origin.

# 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.

$y_i = w^T x_i + w_0$

A better y-intercept
$w_0$

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".

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

$$\underbrace{\qquad}_{X}$$

- 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

Prediction:

$$y = X * w$$

$\underbrace{\qquad}_{\text{Why?}}$

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

$\underbrace{\qquad}_{X}$

# Linear Least Squares

- We can rewrite the objective function as a norm

Want 'w' that minimizes

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

$r_i$

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

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

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

$Xw$

# Linear Least Squares

Want 'w' that minimizes

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

$$\underbrace{\phantom{\frac{1}{2}\|Xw-y\|^2}}_{}$$ Let's expand then compute gradient.

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

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

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

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

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

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

# Linear Least Squares

Training:  `w = solve(X.T @ X,  x.T @ y)`

Why?

Want 'w' that minimizes

$$f(w) = \frac{1}{2}\sum_{i=1}^{n}(w^T x_i - y_i)^2 = \frac{1}{2}\|Xw - y\|_2^2 = \frac{1}{2}w^T X^T X w - w^T X^T y + \frac{1}{2}y^T y \quad \Big\} \text{ a quadratic function (in matrix notation)}$$

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

So at a <u>minimizer</u> where $\nabla f(w) = 0$
we have: $\boxed{X^T X w = X^T y}$     "normal equations"

Some matrix   Some vector

This is a linear system $Aw = b$
for some matrix 'A' and vector 'b'

What are the <u>gradients</u> of these terms?

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

$\nabla_w[w^T b] = b$  →  This is like saying $\frac{d}{dw}[\alpha w] = \alpha$

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

11

# The Punch Line

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

```
w = solve(X.T @ X,  X.T @ y)
```

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

# 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 <u>solutions to the linear system</u>

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

The following are <u>not</u> the solutions to the least squares problem:

$w = (X^T X)^{-1}(X^T y)$  (only true if $X^T X$ <u>is invertible</u>)

$w X^T X = X^T y$  (matrix multiplication is <u>not</u> commutative, dimensions don't even match)

$w = \dfrac{X^T y}{X^T X}$  (you cannot <u>divide by a matrix</u>)

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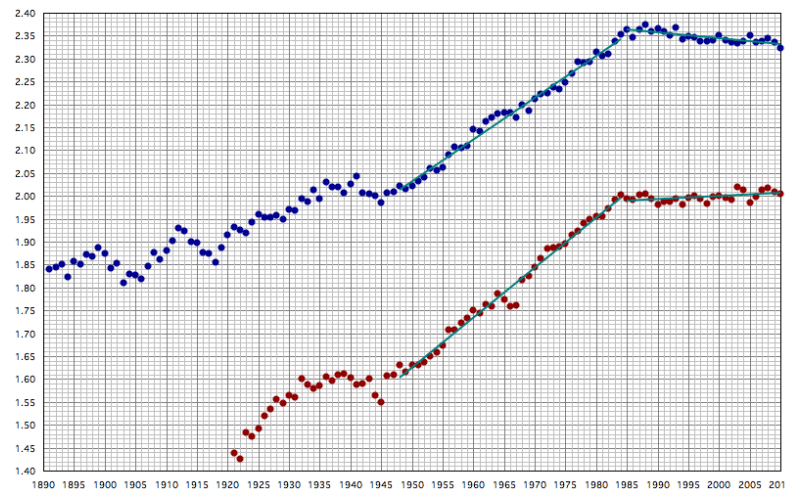
- Each of the $O(d^2)$ elements is an inner product between length 'n' vectors.

# 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)



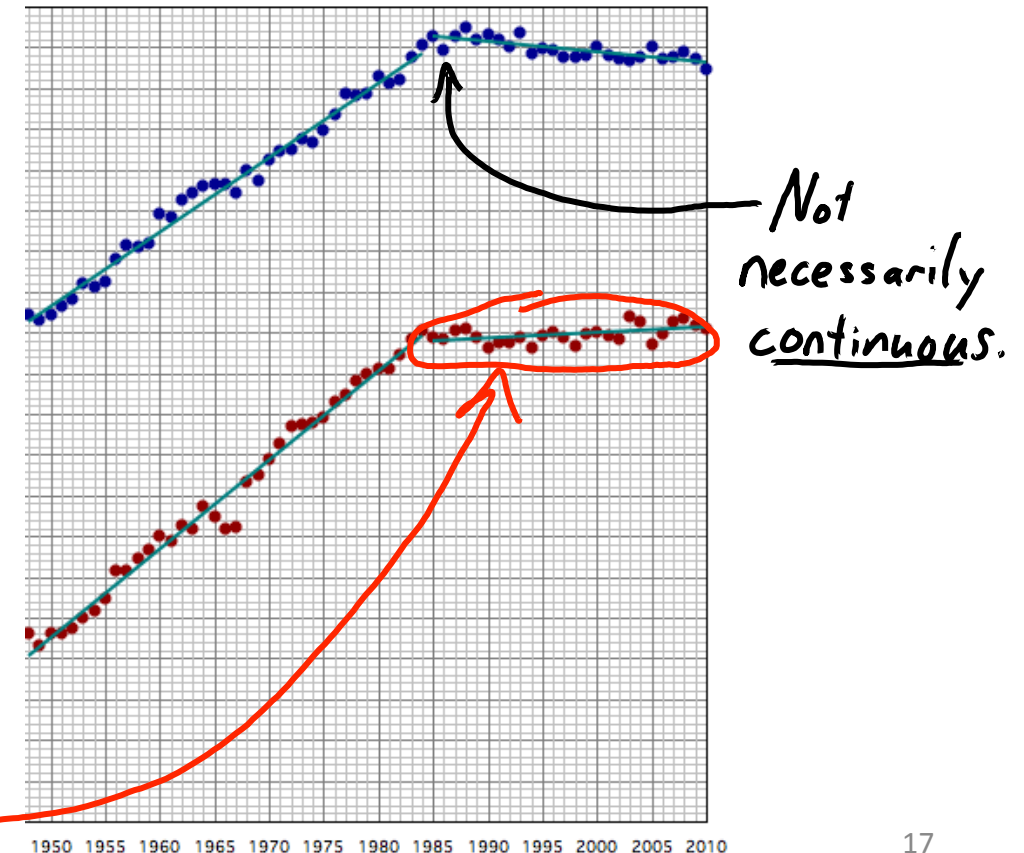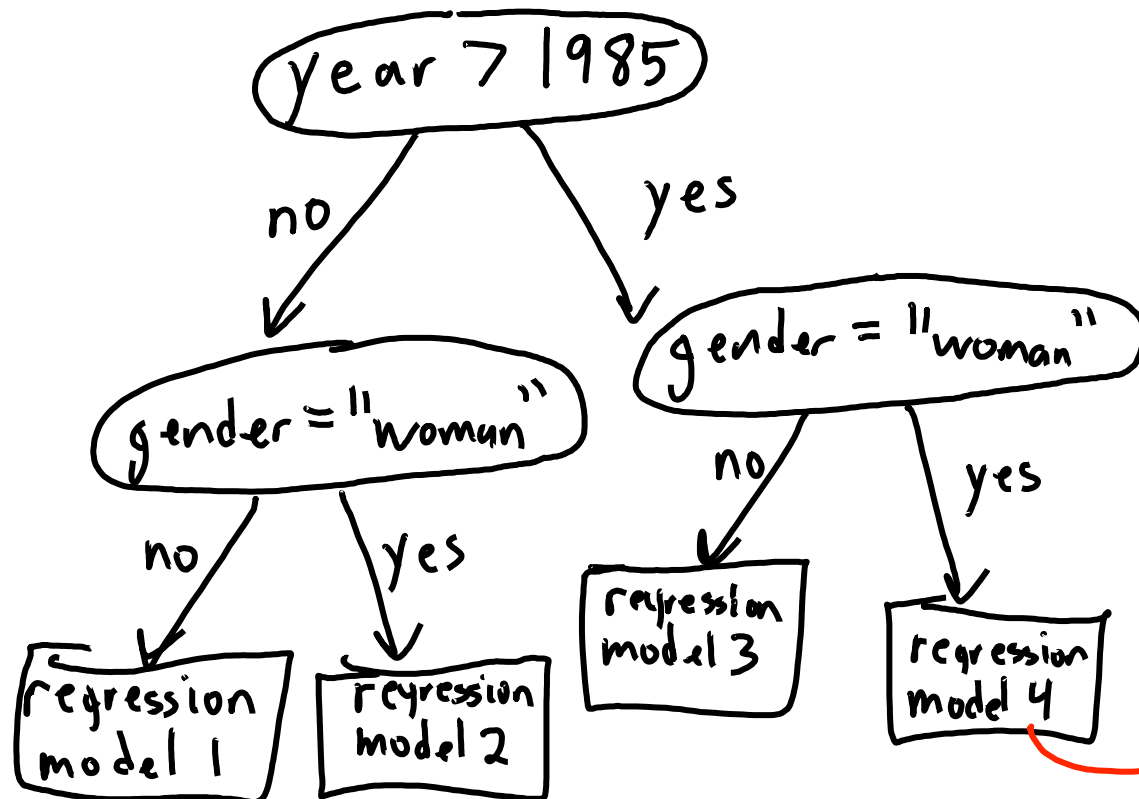SHOT PUT PROGRESSION MEN (7.26 kg) AND WOMEN (4 kg) (mean of top ten)

# Adapting Counting/Distance-Based Methods

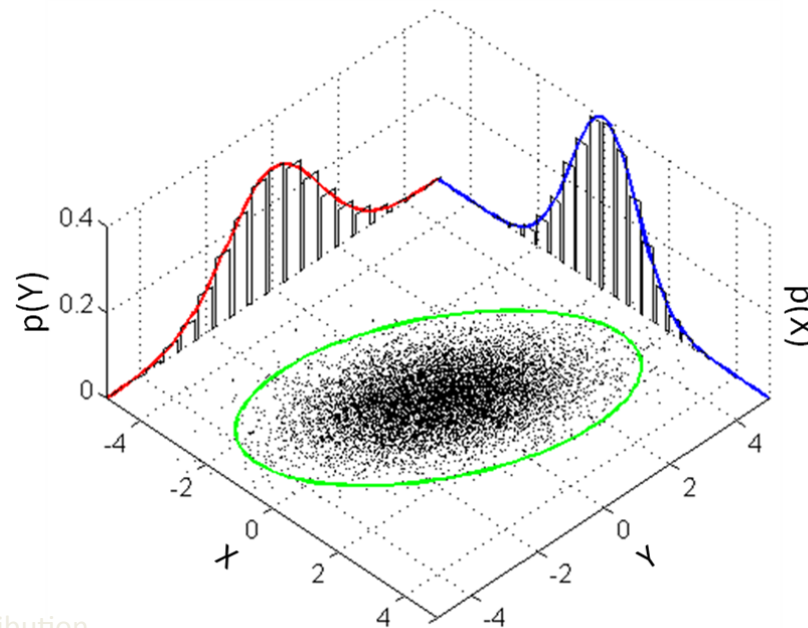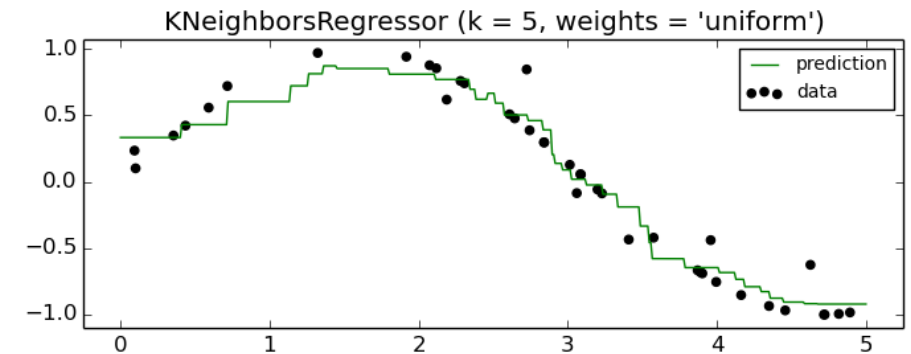- We can adapt our classification methods to perform regression:

# 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.
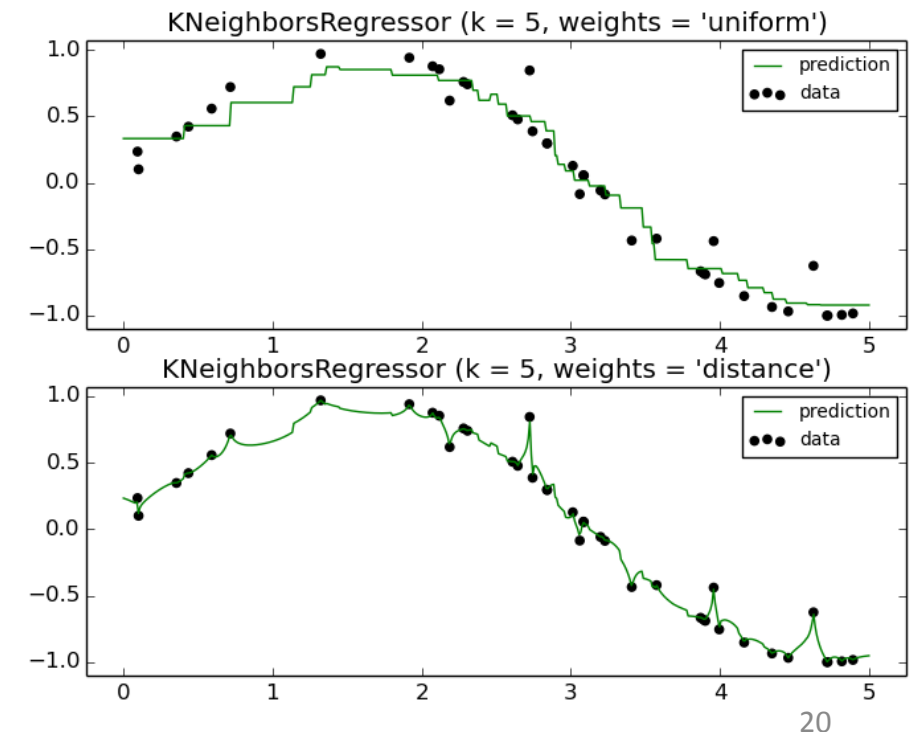
18

# 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')

19

# 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.
  - 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):

$$X = \begin{bmatrix} 0.2 \\ -0.5 \\ 1 \\ 4 \end{bmatrix} \qquad Z = \begin{bmatrix} 1 & 0.2 & (0.2)^2 \\ 1 & -0.5 & (-0.5)^2 \\ 1 & 1 & (1)^2 \\ 1 & 4 & (4)^2 \end{bmatrix}$$
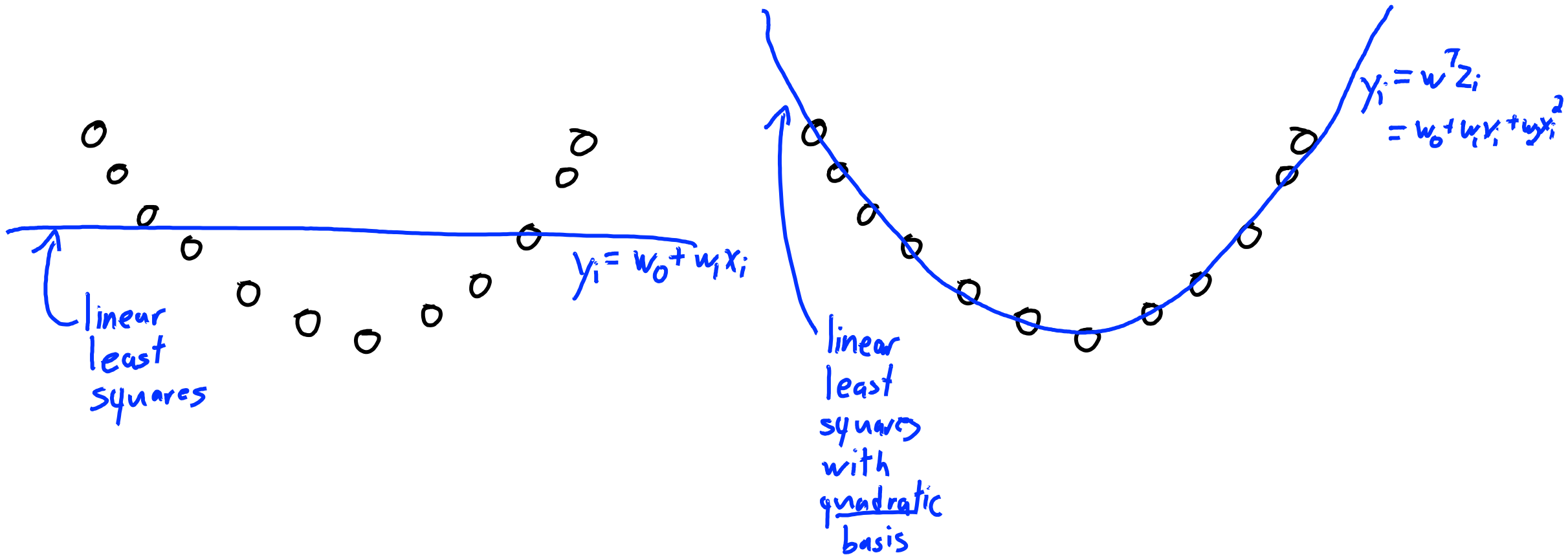
y-int    X    $x^2$

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

"solve linear system"

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

- 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 $\hat{y} = \hat{Z} w$

# Linear Least Squares for Quadratic Models



linear
least
squares

$y_i = w_0 + w_1 x_i$

linear
least
squares
with
quadratic
basis

$y_i = w^T z_i$
$= w_0 + w_1 x_i + w_2 x_i^2$
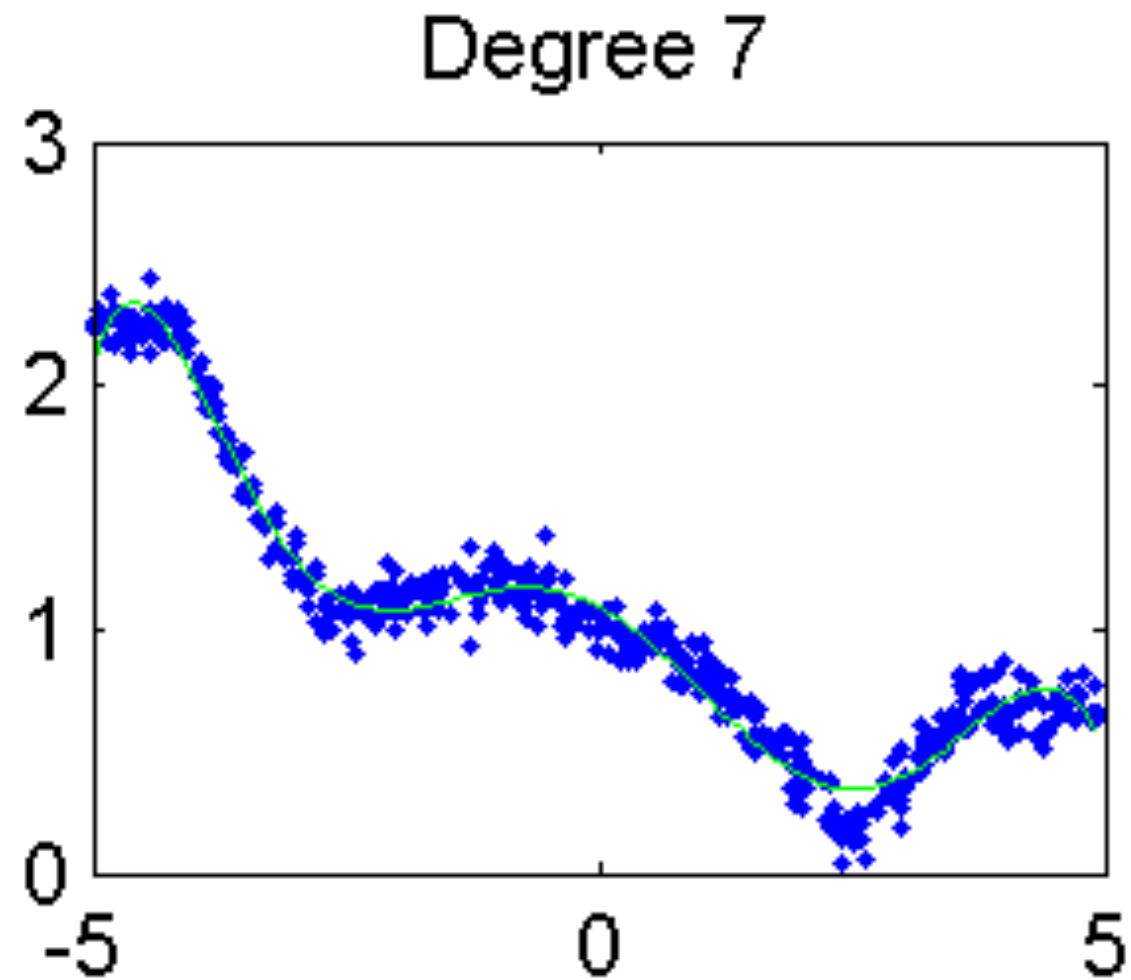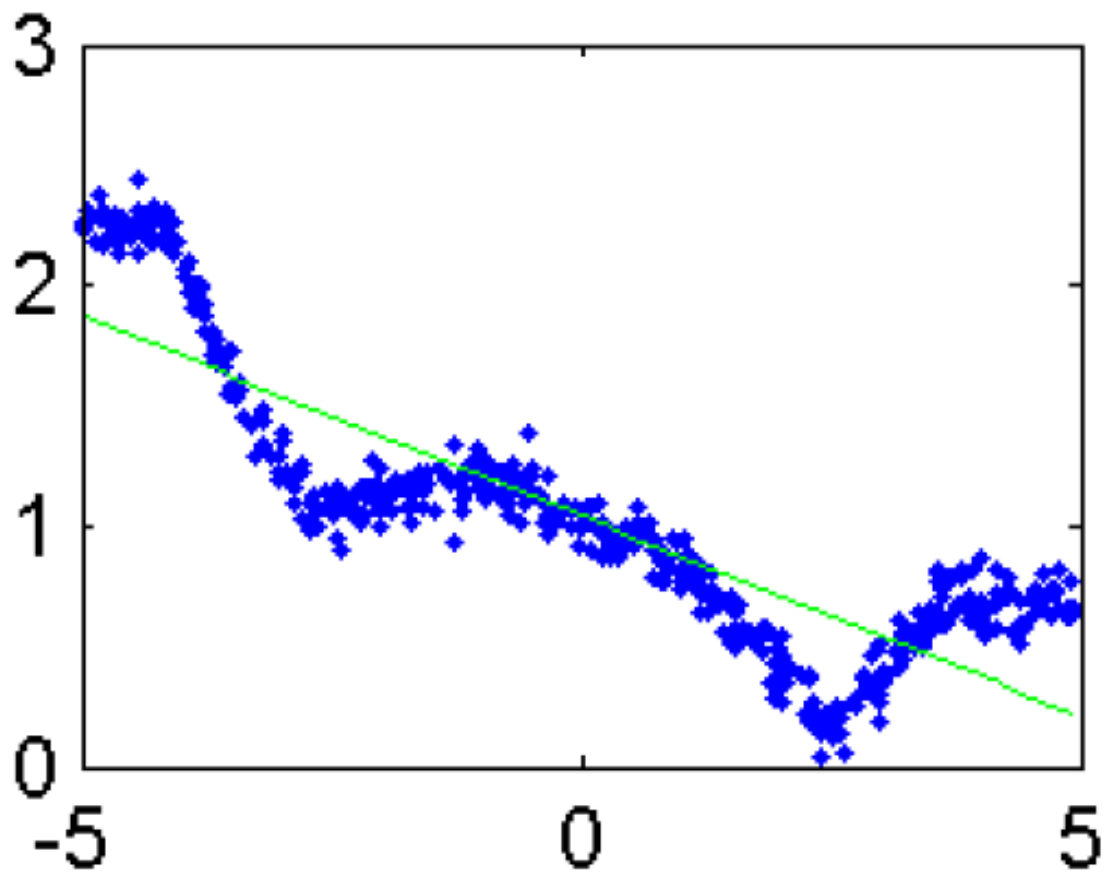
# General Polynomial Basis

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

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

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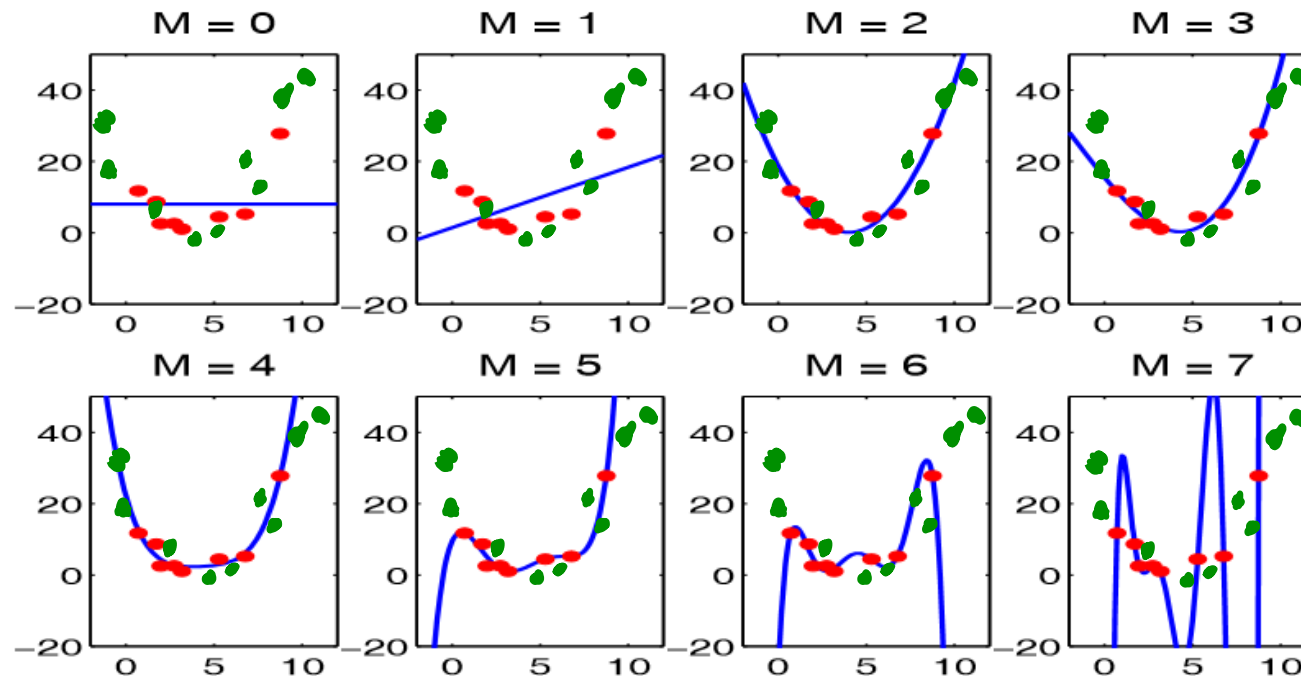  - E.g., Lagrange polynomials (see CPSC 303)

# 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 $(X^TX)w = X^Ty$.
- Tree/generative/non-parametric/ensemble methods for regression.
- Change of basis allows linear models to model non-linear data

- Next time:
  - A general method for avoiding overfitting.