

# Lecture 18: Least squares continued

Thursday, October 29, 2015 9:30 AM

Admin:

## "LEAST SQUARES" FITTING & APPLICATIONS TO DATA ANALYSIS

(Reading: Meyer 4.6  
5.13-14  
Strang 3.3)

Problem: System of equations

$$A\vec{x} = \vec{b}$$

but  $\vec{b} \notin R(A)$ !

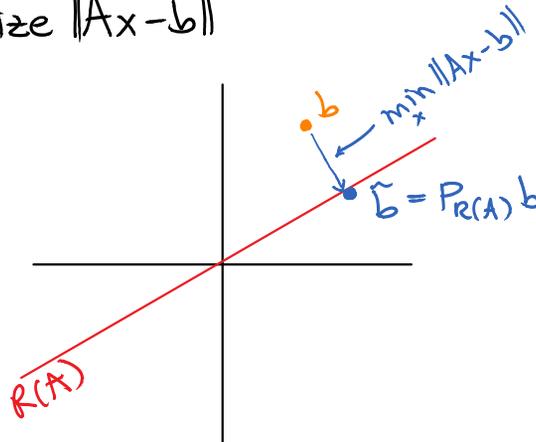
$\Rightarrow$  Find  $x$  to minimize  $\|Ax - b\|$

Algebraically:

$$A = \sum_i \lambda_i \vec{u}_i \vec{v}_i^T$$

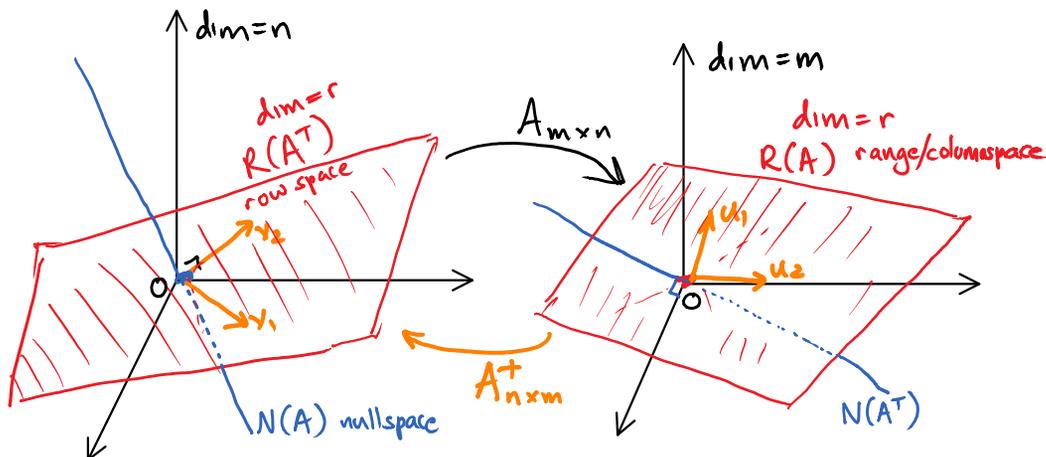
$$A^+ = \sum_{\lambda_i > 0} \frac{1}{\lambda_i} \vec{v}_i \vec{u}_i^T$$

(pseudo-inverse)



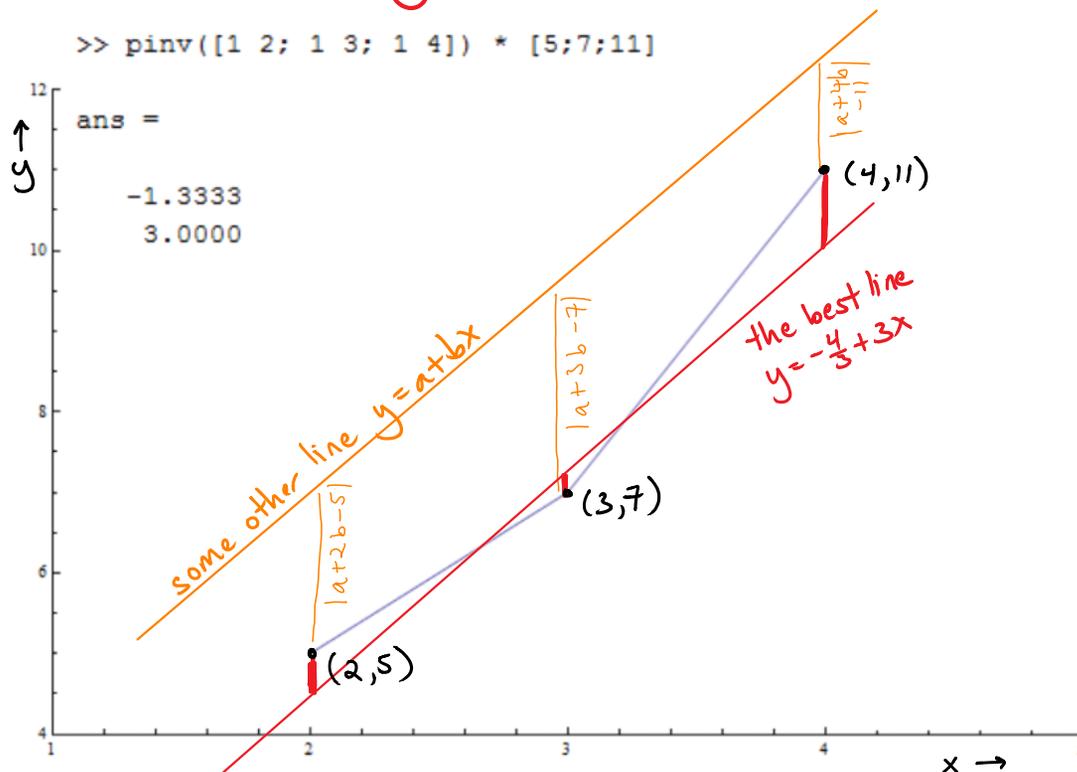
$$\Rightarrow x = A^+ b$$

minimizes  $\|Ax - b\|$



Example: Linear regression

# Example: Linear regression



Setting:  $m$  data points

$$(x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_k^{(1)}, y^{(1)})$$

$$(x_1^{(2)}, x_2^{(2)}, \dots, x_k^{(2)}, y^{(2)})$$

$$\vdots$$

$$(x_1^{(m)}, x_2^{(m)}, \dots, x_k^{(m)}, y^{(m)})$$

components known exactly, eg., date/time

component that we want to predict

Goal: Find the best linear predictor for  $y$ ,

$$a_0, a_1, a_2, \dots, a_k \in \mathbb{R}$$

to minimize total squared error.

$$\sum_{j=1}^m |a_0 + a_1 x_1^{(j)} + \dots + a_k x_k^{(j)} - y^{(j)}|^2$$

Answer:

$$\begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_k^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_k^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_k^{(m)} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \end{pmatrix} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix}$$

$$\begin{pmatrix} 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_k^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} & \dots & x_k^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_k^{(m)} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{pmatrix} = \begin{pmatrix} y^{(2)} \\ y^{(3)} \\ \vdots \\ y^{(m)} \end{pmatrix}$$

$A$   $y$

The least-squares solution is  $A^+ \vec{y}$ .  
↑  
pseudoinverse

### Four ways to find $x$ to minimize $\|Ax - b\|$ :

① Compute pseudo-inverse  $A^+$ , set  $x = A^+ b$ .

**Pros:** Easy to remember, one line `pinv(A)*b` in Matlab

**Cons:** Slow for large  $A$ , numerically unstable

② Find  $x$  manually, with calculus:

$$E = \left\| \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_m \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} \right\|^2$$

squared error

$$= \sum_{j=1}^m (x_1 + a_j x_2 - b_j)^2$$

$$\frac{\partial E}{\partial x_1} = 2 \sum_j (x_1 + a_j x_2 - b_j) = 0 \Rightarrow m x_1 + (\sum_j a_j) x_2 = (\sum_j b_j)$$

$$\frac{\partial E}{\partial x_2} = 2 \sum_j a_j (x_1 + a_j x_2 - b_j) = 0 \Rightarrow (\sum_j a_j) x_1 + (\sum_j a_j^2) x_2 = \sum_j a_j b_j$$

Let  $\bar{a} = \frac{1}{m} \sum_j a_j$  (average value of  $a_j$ )  
 $\bar{b} = \frac{1}{m} \sum_j b_j$  (average of  $b_j$ 's)

First equation

$$\Rightarrow \text{intercept } x_1 = \bar{b} - \bar{a} x_2$$

(this is 0 if the data is centered so  $\bar{a} = \bar{b} = 0$ )

Second equation

$$\Rightarrow x_2 = \frac{\sum_j a_j b_j - m \cdot \bar{a} \cdot \bar{b}}{\sum a_j^2 - m(\bar{a})^2}$$

$$\left( = \frac{\text{Cov}(A, B)}{\text{Var}(A)} \quad \text{if } A \text{ and } B \text{ are random variables} \right. \\ \left. \text{with } \Pr[A, B] = (a_j, b_j) = \frac{1}{m}. \right)$$

**Pros:** None (okay, it gives an easy closed-form solution)

③ Solve  $A^T A x = A^T b$

Theorem: Let  $A$  be an  $m \times n$  real matrix with rank  $n$ .  
Then for any  $b$ ,

- The equation  $A^T A x = A^T b$  is feasible (ie, it has a solution  $x$ )
- The unique solution is

$$\boxed{(A^T A)^{-1} A^T b = A^+ b}$$

↑  
pseudoinverse

Proof: There is a solution because

- $R(A^T A) = R(A^T)$  (shown above)
- $A^T b \in R(A^T)$

$$\Rightarrow A^T b \in R(A^T A) \quad \checkmark$$

The solution is unique because  $N(A) = \{\vec{0}\}$ .  
(because  $\text{rank}(A) = \dim R(A^T) = n$ , by assumption,  
and by rank-nullity  $\dim R(A^T) + \dim N(A) = n$ )

So why is  $(A^T A)^{-1} A^T b = A^+ b$ , as claimed?  
There are several ways of seeing it.

Algebraically: Using the SVD of  $A$ ,  $A = \sum_i \lambda_i \vec{u}_i \vec{v}_i^T$ ,

$$A^+ = \sum_{i: \lambda_i > 0} \frac{1}{\lambda_i} \vec{v}_i \vec{u}_i^T$$

$$\begin{aligned} (A^T A)^{-1} A^T &= \left[ \sum_i \lambda_i^2 \vec{v}_i \vec{v}_i^T \right]^{-1} \left( \sum_j \lambda_j \vec{v}_j \vec{u}_j^T \right) \\ &= \left( \sum_i \frac{1}{\lambda_i^2} \vec{v}_i \vec{v}_i^T \right) \left( \sum_j \lambda_j \vec{v}_j \vec{u}_j^T \right) \end{aligned}$$

↑ since  $\text{rank}(A) = n$ ,  
↑ 1 - - - - - 11 - - - - - n

$$\begin{aligned}
 & \text{since rank}(A) = n, \\
 & \lambda_1, \dots, \lambda_n \text{ are all } > 0 \\
 & = \sum_i \frac{1}{\lambda_i} v_i u_i^T \\
 & = A^+ \quad \checkmark
 \end{aligned}$$

More geometrically:

$$\begin{aligned}
 x = A^+ b \text{ solves } Ax &= P_{R(A)} b \\
 \Rightarrow A^T A x &= A^T P_{R(A)} b \\
 &= A^T b \quad \text{since } A^T = A^T P_{R(A)} \\
 & \text{(in general } A = P_{R(A)} A P_{R(A^T)})
 \end{aligned}$$

Corollary: For any  $m \times n$  real matrix  $A$ , of rank  $n$ ,

$$A^+ = (A^T A)^{-1} A^T$$

Proof: Since we just showed  $A^+ b = (A^T A)^{-1} A^T b$  for any  $b$ .  $\square$

Pros: Fast

Cons: condition number of  $A^T A$   
 $= (\text{condition number of } A)^2$

Method  
 ④ Solve  $\begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} y \\ x \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix}$

Pros: Maintains sparsity of  $A$ .

[Problem 4.6.9, p.237 of Meyer]

⑤ Latest research has focused on finding fast randomized approximation algorithms, based on dimension reduction à la Johnson-Lindenstrauss Lemma, eg.,

[Clarkson Woodruff 2012 <http://arxiv.org/abs/1207.6365>]

[Nelson & Nguyen 2012 <http://arxiv.org/abs/1211.1002>]

## Evaluating the quality of a fit:

Of course there is a lot of statistics to consider when running

Of course there is a lot of statistics to consider when running linear regressions.

The most basic statistic is the " $R^2$  value", also known as the "coefficient of determination."

Definition: Given data points  $(x_1, y_1), \dots, (x_m, y_m)$ , and the least-squared-error fit line  $y = \alpha + \beta x$ ,

$$R^2 = 1 - \frac{\sum_j (y_j - \alpha - \beta x_j)^2}{\sum_j (y_j - \bar{y})^2}$$

$$= 1 - \frac{\text{(total squared error from best fit line)}}{\text{(total squared error from horizontal line } \bar{y} \text{)}} \quad \text{where } \bar{y} = \frac{1}{m} \sum_j y_j$$

Observe:  $0 \leq R^2 \leq 1$

and higher is better

- Collinear data points will have  $R^2 = 1$  (since best-fit line will fit exactly)
- But what counts as a "good fit" depends very much on the field (eg., macroeconomics)

Interpretation of  $R^2$ : Explained variance

If the data were distributed randomly about the mean  $\bar{y}$ , the observed variance is  $\mathbb{E}[Y - \mathbb{E}(Y)]^2 = \frac{1}{m} \sum_{j=1}^m (y_j - \bar{y})^2$ .

For data distributed randomly off the line  $y = \alpha + \beta x$ ,

ie.,  $y_j = \alpha + \beta x_j + e_j$ , the remaining variance is  $\frac{1}{m} \sum_j e_j^2 = \frac{1}{m} \sum_j (y_j - \alpha - \beta x_j)^2$ .

$\Rightarrow R^2$  measures "how much of the observed variance is explained by the observed  $x_j$  data points."

### Probabilistic interpretation: Why least squares?

Choosing  $\vec{\alpha}$  to minimize  $\|A\vec{\alpha} - \vec{y}\|$  is geometrically natural, but why is this natural in data analysis?

One reason (of several): In a natural probabilistic model,

least-squares regression finds the **maximum likelihood** estimate of  $\vec{a}$ .

Model: Assume that the  $y$ -values are generated via

$$y^{(i)} = \vec{a} \cdot \vec{x}^{(i)} + \varepsilon^{(i)}$$

where  $\varepsilon^{(i)} \sim N(0, \sigma^2)$  — independent, identically distributed, Gaussian random noise.

$$\Rightarrow p_{\vec{a}}(y^{(i)} | x^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(y^{(i)} - \vec{a} \cdot x^{(i)})^2\right]$$

↑ probability density function of  $y^{(i)}$   
for model parameters  $\vec{a}$

$$\Rightarrow p_{\vec{a}}(y^{(1)}, \dots, y^{(m)} | x^{(1)}, \dots, x^{(m)}) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \vec{a} \cdot x^{(i)})^2\right]$$

⇒ Setting  $\vec{a}$  to minimize  $\sum_i (y^{(i)} - \vec{a} \cdot x^{(i)})^2$   
**maximizes the likelihood** of the observed data.

(Notice: This holds for any value of the variance  $\sigma^2$ ,  
and  $\sigma$  need not be known to choose  $\vec{a}$ .)

### Linear regression with weights:

If errors between observations are correlated, or if some observations are more reliable than others, then minimizing the weighted sum of squared errors is better than unweighted.

Example: Say

$$y_1 = \alpha x_1 + \varepsilon_1, \quad \varepsilon_1 \sim N(0, 1) \quad \leftarrow \text{more reliable}$$

$$y_2 = \alpha x_2 + \varepsilon_2, \quad \varepsilon_2 \sim N(0, 2) \quad \leftarrow \text{less reliable}$$

$$\Rightarrow p_{\alpha}[y_1, y_2 | x_1, x_2] = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y_1 - \alpha x_1)^2} \cdot \frac{1}{\sqrt{2\pi \cdot 2}} e^{-\frac{1}{2 \cdot 2}(y_2 - \alpha x_2)^2}$$

$$\propto \exp\left[-\frac{1}{2} \cdot \left\{ (y_1 - \alpha x_1)^2 + \frac{1}{2} (y_2 - \alpha x_2)^2 \right\}\right]$$

⇒ The max-likelihood  $\alpha$  minimizes the weighted sum  
 $S = \frac{1}{\sigma_1^2} (y_1 - \alpha x_1)^2 + \frac{1}{\sigma_2^2} (y_2 - \alpha x_2)^2$

$$\frac{-1}{2} \frac{\partial S}{\partial \alpha} = \frac{x_1}{\sigma_1^2} (y_1 - \alpha x_1) + \frac{x_2}{\sigma_2^2} (y_2 - \alpha x_2) = 0$$

$$\Rightarrow \alpha = \frac{\frac{x_1 y_1}{\sigma_1^2} + \frac{x_2 y_2}{\sigma_2^2}}{\frac{x_1^2}{\sigma_1^2} + \frac{x_2^2}{\sigma_2^2}}$$

More generally: To choose  $\alpha$  to minimize

$$\sum_{j=1}^m w_j (y_j - x_j \alpha)^2,$$

observe

$$\| W(y - A\alpha) \|^2$$

$$\text{where } W = \begin{pmatrix} w_1 & & 0 \\ & w_2 & \\ 0 & & w_m \end{pmatrix}$$

$$A = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$$

This is standard least-squares for the equation

$$WA\alpha = Wy$$

$\Rightarrow$  best  $\alpha$  satisfies

$$A^T W^T W A \alpha = A^T W^T W y, \text{ etc.}$$

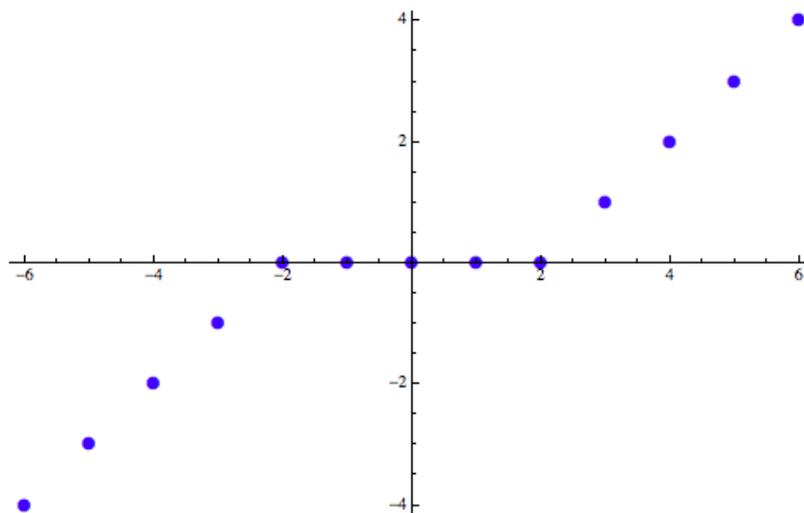
Practical application: **Locally weighted linear regression**

Example:

```
data = {
  {-6, -4},
  {-5, -3},
  {-4, -2},
  {-3, -1},
  {-2, 0},
  {-1, 0},
  {0, 0},
  {1, 0},
  {2, 0},
  {3, 1},
  {4, 2},
  {5, 3},
  {6, 4}
};
```

```
m = Length[data];
```

```
dataplot = ListPlot[data]
```



```
X = Table[{1, data[[j, 1]]}, {j, 1, m}];
```

```
X
```

```
y = data[[All, 2];
```

```
ymean =  $\frac{1}{m}$  Plus@@y;
```

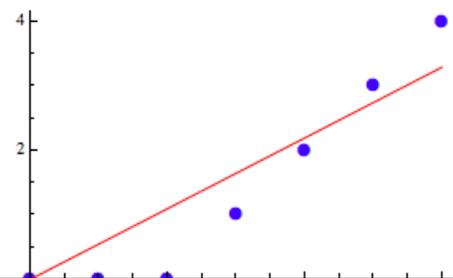
```
linear = PseudoInverse[X].y
```

```
plotlinear = Plot[{1, x}.linear, {x, -6, 6}
```

```
Show[dataplot, plotlinear]
```

```
{{1, -6}, {1, -5}, {1, -4}, {1, -3}, {1, -2}, {1, -1}, {1, 0}, {1, 1}, {1, 2}, {1, 3}, {1, 4}, {1, 5}, {1, 6}}
```

**Linear regression**



```
linear = PseudoInverse[a];
plotlinear = Plot[{1, x}.linear, {x, -6, 6}
Show[dataplot, plotlinear]
```

```
{(1, -6), (1, -5), (1, -4), (1, -3), (1, -2),
(1, -1), (1, 0), (1, 1), (1, 2), (1, 3), (1, 4), (1, 5), (1, 6)}
```

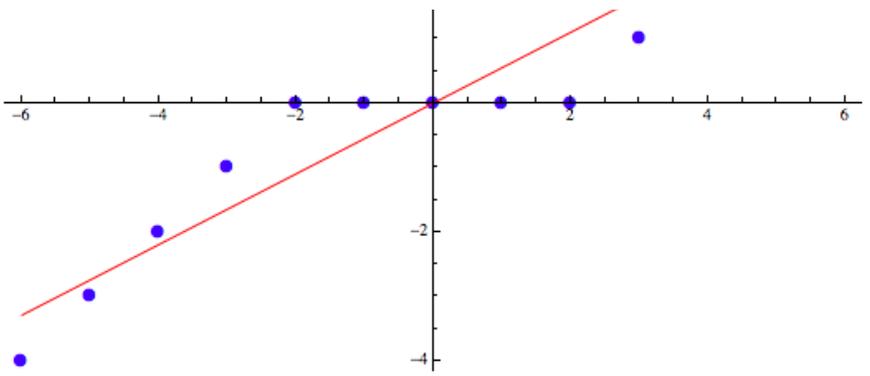
```
{0, 50/91}
```

"R^2 value:"

$$R2linear = 1 - \frac{\sum_{j=1}^n (y[j] - X[j].linear)^2}{\sum_{j=1}^n (y[j] - ymean)^2} // 1$$

R^2 value:

0.915751



Cubic regression  $1, x, x^2, x^3$

```
HigherOrderRegression[data, 3]
```

```
α = {0, 0.1998, 0, 0.013986}
```

```
R2 = 0.982884
```

```
HigherOrderRegression[data_, order_] :=
```

```
Module[{m, y, ymean, X, α, plot, min, max, R2},
```

```
m = Length[data];
```

```
y = data[All, 2];
```

```
X = Table[Prepend[Table[data[[j, 1]]^k, {k, 1, order}], 1], {j, 1, m}] // N;
```

```
α = PseudoInverse[X].y;
```

```
Print["α = ", α // Chop];
```

```
{min, max} = {Min[data[All, 1]] - 2, Max[data[All, 1]] + 2};
```

```
plot = Plot[Table[x^k, {k, 0, order}].α, {x, min, max}];
```

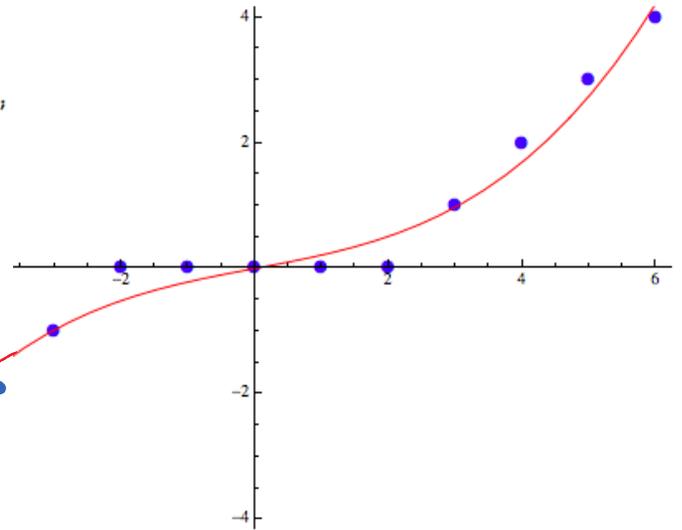
```
ymean = 1/m Plus@@y;
```

$$R2 = 1 - \frac{\sum_{j=1}^n (y[j] - X[j].α)^2}{\sum_{j=1}^n (y[j] - ymean)^2} // N;$$

```
Print["R2 = ", R2];
```

```
Show[dataplot, plot]
```

```
];
```

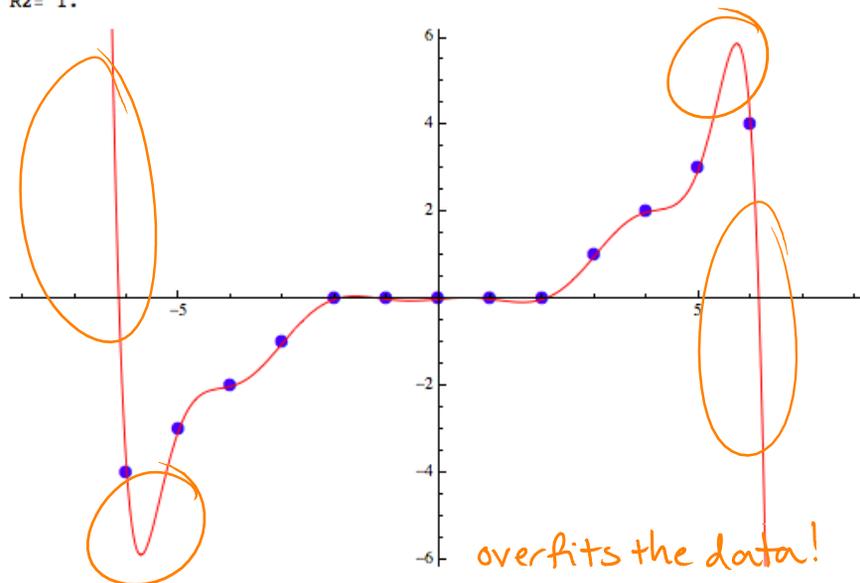


11<sup>th</sup> order  $(1, x, x^2, \dots, x^{11})$

```
Show[HigherOrderRegression[data, 11], PlotRange -> {{-8, 8}, {-6, 6}}]
```

```
α = {1.70782 × 10-10, 0.101443, -2.42844 × 10-10,
-0.137535, 0, 0.03917, 0, -0.00318122, 0, 0.000104718, 0, -1.2025 × 10-6}
```

R2 = 1.



overfits the data!

Locally weighted linear regression

— to predict  $y$  at point  $x$ ,  
 run weighted linear regression  
 giving higher weights to data points close to  $x$ !  
 e.g.,  $w_j = e^{-\frac{1}{2\tau^2}(x_j - x)^2}$

```

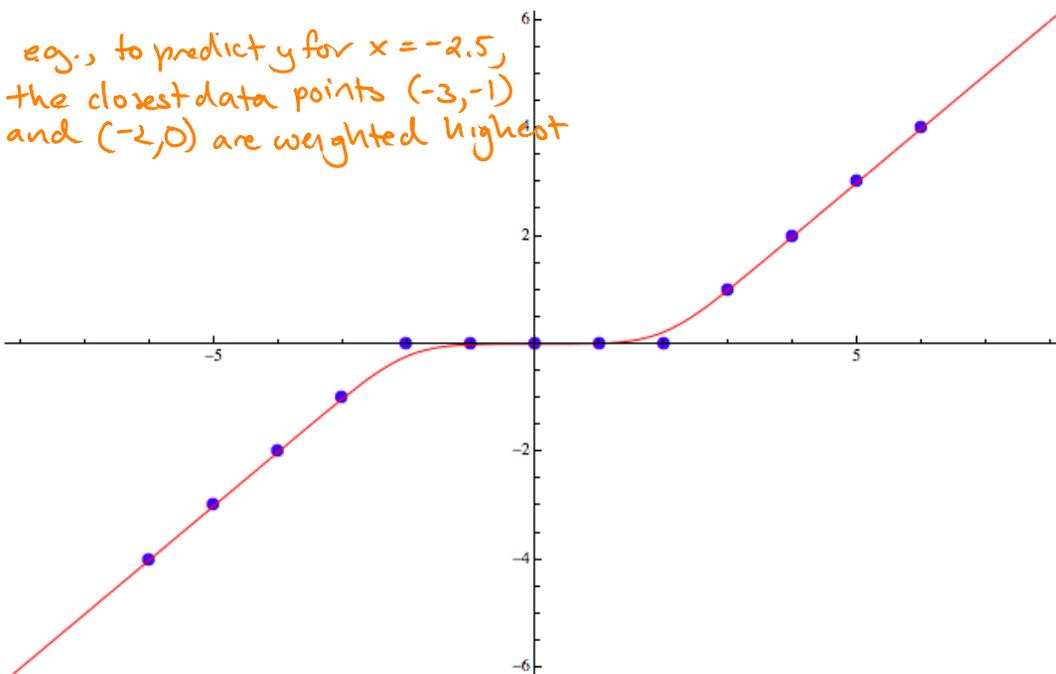
m = Length[data];
y = data[[All, 2]];
X = Table[{1, data[[j, 1]]}, {j, 1, m}] // N[#, 32] &;

τ = 1;

weightedlinearplot = Plot[
  W = DiagonalMatrix[Table[Exp[-1/(2 τ^2) (data[[j, 1]] - x)^2], {j, 1, m}]];
  α = PseudoInverse[W.X].W.y;
  {1, x}.α,
  {x, -10, 10}];
Show[dataplot, weightedlinearplot, PlotRange -> {{-8, 8}, {-6, 6}}]

```

e.g., to predict  $y$  for  $x = -2.5$ ,  
 the closest data points  $(-3, -1)$   
 and  $(-2, 0)$  are weighted highest



Note: • After running weighted linear regression at  $x$ ,  
 we don't keep the line — just the prediction at  $x$ .

- This is slower than running linear regression just once.
- Full "training set" must be kept to make predictions.  
 (known as a "non-parametric" learning algorithm)