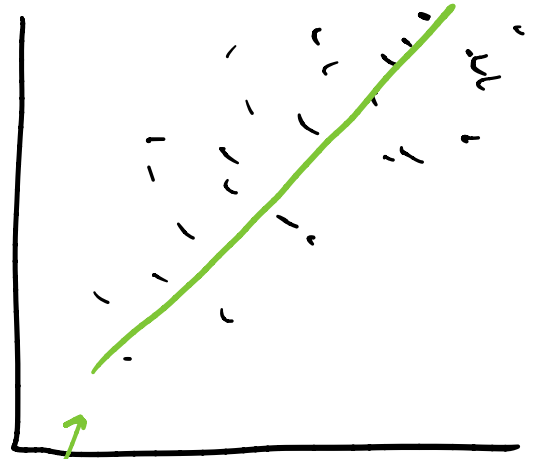
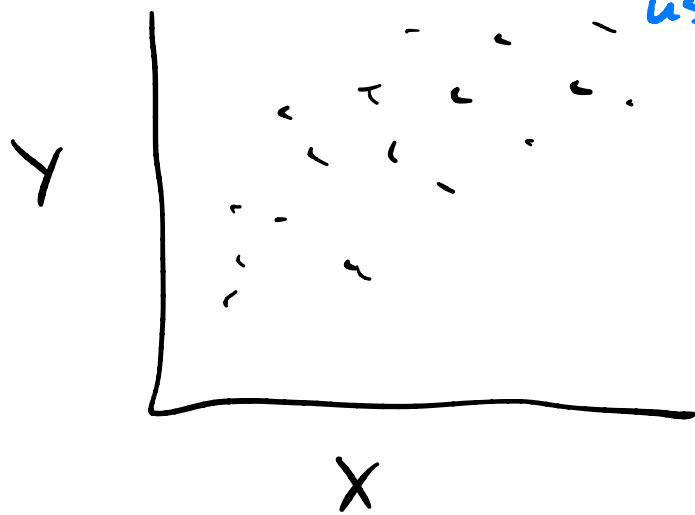


Lecture 12: Linear Regression

① $F(x_i) = b + m x_i$ ← know this — data

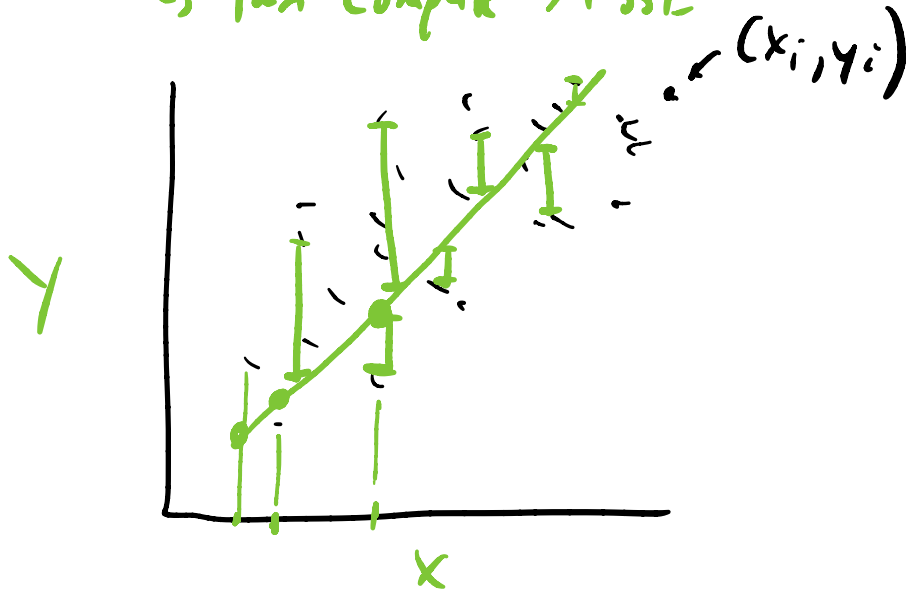
← solve for these unknowns using SciPy



② Compute MSSE, regression line

Option 1: plug in M x_i data used to find b & m into $F(x_i) = b + m x_i$ ← data used to fit regression

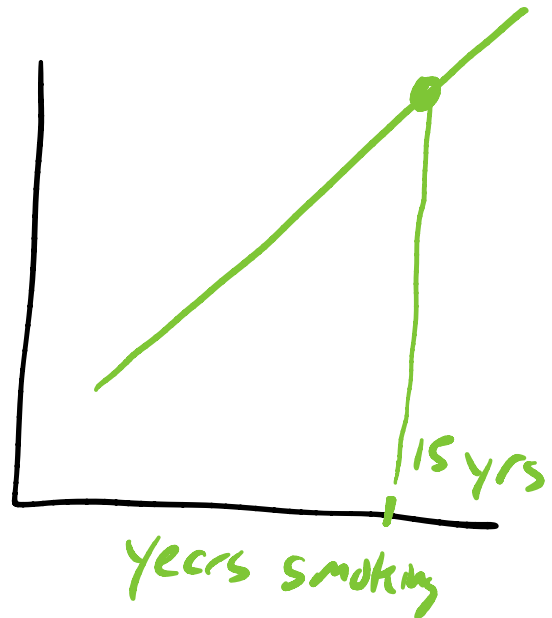
↳ then compute MSSE



Option 2: predict data values we've never seen before



Cancer
rate



$$F(15) = b + m \cdot 15 \rightarrow \text{prediction}$$

$$\begin{array}{c} \overbrace{F(x_i)}^{\vec{y}} \\ F(x_i) = b + m \overbrace{x_i}^{\text{ind var}} \end{array} \quad \left. \begin{array}{c} \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{array} \right] \right. \\ \underbrace{\hspace{1.5cm}}_A \end{array} \right] \text{one sample}$$

$$\left[\begin{array}{c} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{array} \right]$$

$$\vec{y} = A \vec{c}$$

$b \cdot (1) + m \cdot x_i$
 "homogeneous coord" — col of 1s
 used for the intercept $b \rightarrow$ applies it to x_i value.

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \begin{bmatrix} b \\ m \end{bmatrix}$$

$\vec{c}: (2, 1)$
 $A: (N, 2)$
 Shape: $(N, 1)$

$1 \cdot b + m \cdot x_1 = y_1$
 $1 \cdot b + m \cdot x_2 = y_2$

Rename :

$$b \rightarrow c_0$$

$$m \rightarrow c_1$$

$$F(x_i) = c_0 + c_1 x_i$$

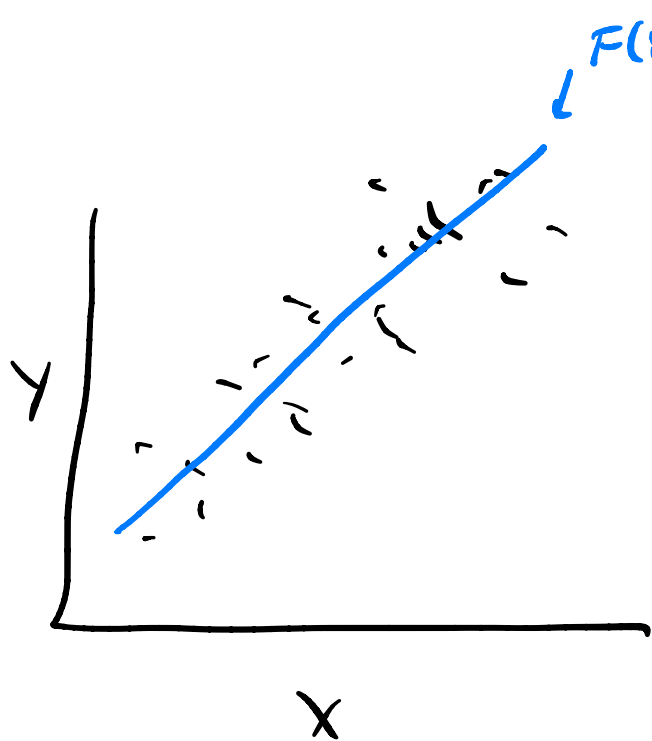
$$\vec{c} = \begin{bmatrix} c_0 \\ c_1 \end{bmatrix}$$

Shape: $(2, 1)$

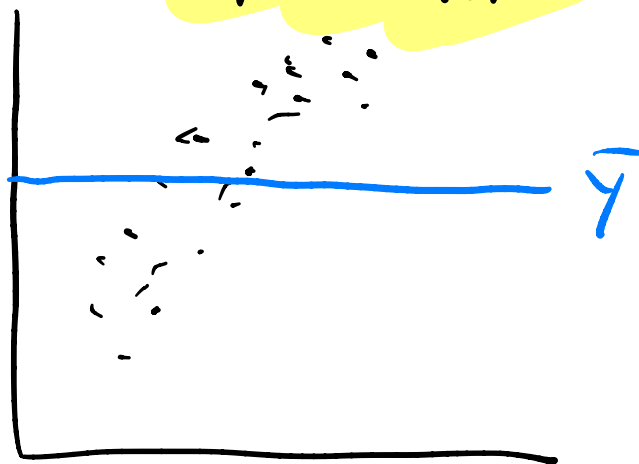
Quality of fit

Is bothering to do a regression worth the effort?

⇒ we compare doing regression with simply calculating the mean of the dependent var Y : \bar{Y}



vs.



$$\bar{Y} = \frac{1}{N} \sum_{i=1}^N Y_i$$

predicted values from regression $\hat{Y}_i = F(x_i) = \underbrace{c_0}_{\text{predicted Value}} + \underbrace{c_1}_{\text{predicted Value}} x_i$

Compare with mean of Y_i values \bar{Y}

How much of improvement \hat{Y}_i vs. \bar{Y}

$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

predicted vals from regression compared to mean

discrepancy of y_i values compared to their mean.

also written:

Sum of squared errors (SSE)

$$R^2 = 1 - \frac{\|r\|_2^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

residuals

$\|r\|_2 = \text{distance}$

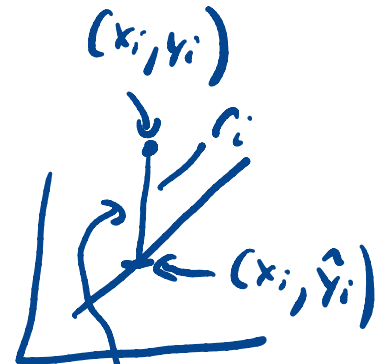
$$= \sqrt{r_1^2 + r_2^2 + r_3^2 + \dots + r_N^2}$$

$$\|r\|_2^2 = r_1^2 + r_2^2 + r_3^2 + \dots + r_N^2$$

error

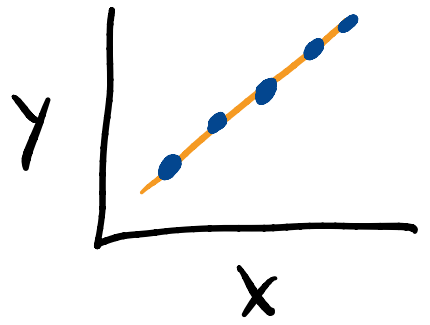
Sum of Squares Error

SSE



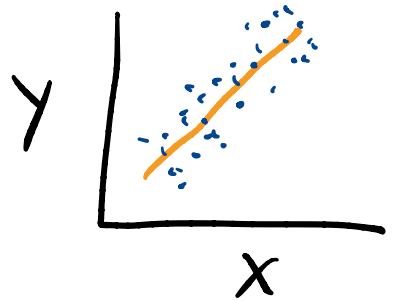
$$1) \overset{sse}{\underbrace{\|r\|_2^2}} \approx 0, R^2 = 1$$

NO error!



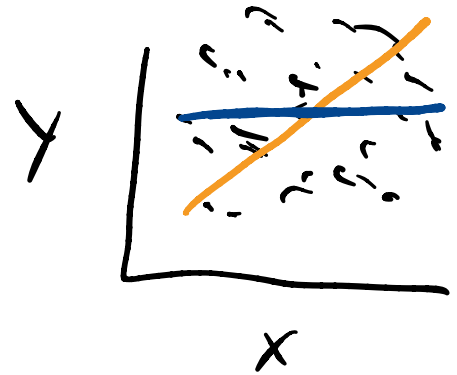
$$2) \left. \begin{array}{l} \text{error is small} \\ \|r\|_2^2 \end{array} \right\} \text{ is } \underline{\text{much}} \underline{\text{less}} \text{ than } \sum_{i=1}^n (y_i - \bar{y})$$

$R^2 \approx 1$



$$3) \|r\|_2^2 \approx \sum_{i=1}^n (y_i - \bar{y})$$

$R^2 \approx 0$



regression not useful/worth it