

Linear Regression

15 December 2023 07:52

model is linear \rightarrow linear relationship between Y & X_i 's

Supervised Learning Algorithm

\rightarrow continuous / categorical



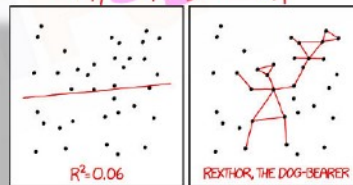
Linear Regression using Python



What is Regression?

- A technique of finding the relationship between two or more variables
- Change in dependent variable is associated with a change in one or more independent variables.

$x_1, x_2, x_3, \dots, x_n$



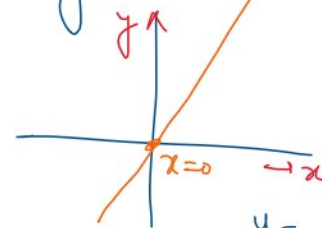
I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

www.intellipaat.com

a function in x

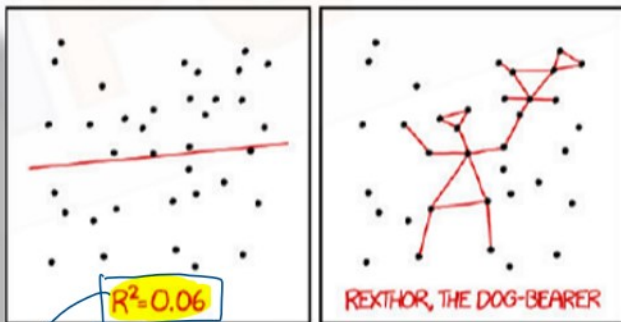
$$y = 2x$$

$$y = mx$$



$$y = 2x$$

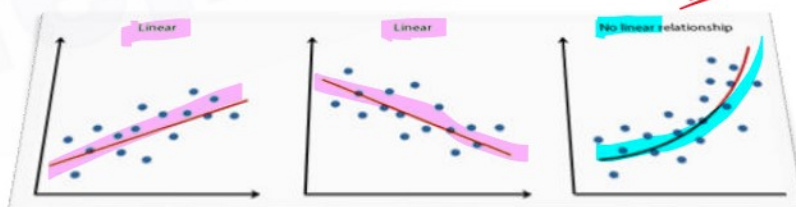
y	x
0	0
-2	-1
2	1



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

6%

Non-linear



function

① Intuitive Approach
↓

② Mathematical Approach
↓

③ Python hands-on

- ①
- ②
- ③
- ④

Calculus

Linear Algebra

Probability and Random Variables
⊕
Probability distribution

statistic

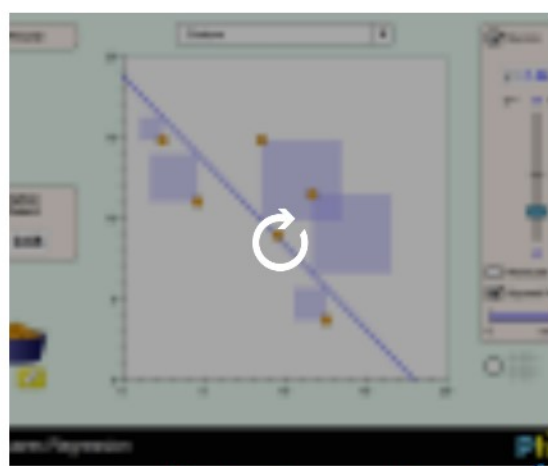
function

limits → Continuity → Differentiability

Maxima & Minima (Global Optimization)

Normal distribution

Least-Squares Regression



mathematical equation

$$Y = 0.83X + 0.75$$

$$Y = \hat{\beta}_1 x + \hat{\beta}_0$$

slope intercept

$Y = 0.83X + 0.75 \Leftrightarrow Y = (mX + c)$: slope - intercept line.

Best-Fit Line

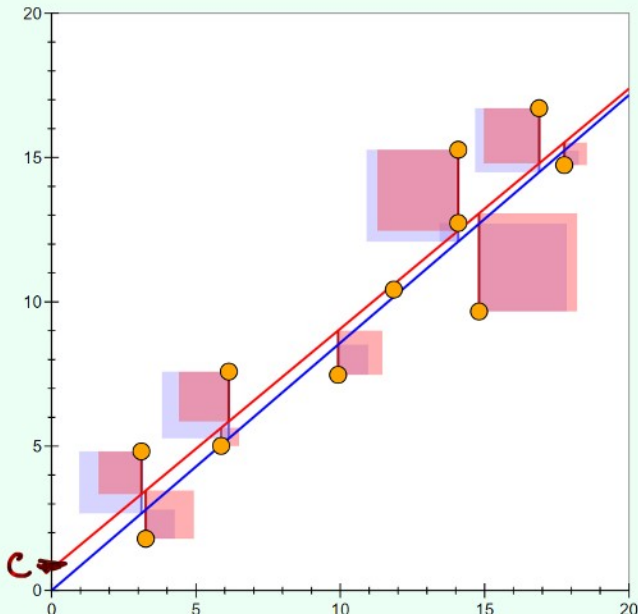
☒ Best-Fit Line

$y = 0.83x + 0.75$

☒ Residuals

☒ Squared Residuals

Correlation Coefficient



☒ My Line

$y = 0.86x + 0.00$

$y = a x + b$

☒ Residuals

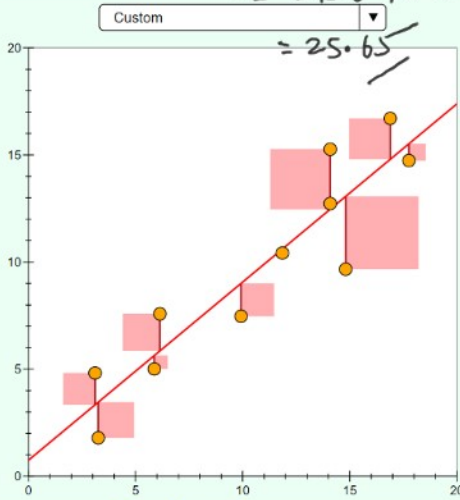
☒ Squared Residuals

$$y = 0.83x + 0.75$$

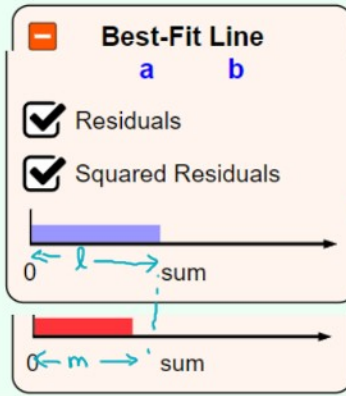
$$y = 0.83 \times 30 + 0.75$$

$$\text{Predicted value} = 24.90 + 0.75$$

at $x = 30$



• : data points
— : Linear Regression Line



l : intuitive best fit line

m : LR model best fit line

a) $l > m$

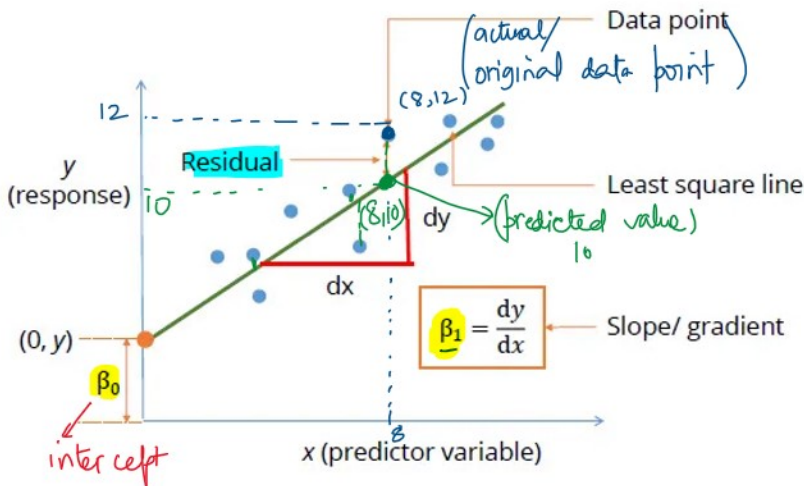
b) $l < m$

c) $l = m$

d) Not sure

Best future

• data point
— (Linear Regression/ Best fit/ least square line)



$$\text{(Predicted - actual)}$$

$$(10 - 12)$$

$$= -2$$

+ve residues cancel -ve residues

$$y = \underline{c} + \underline{m}x + u$$

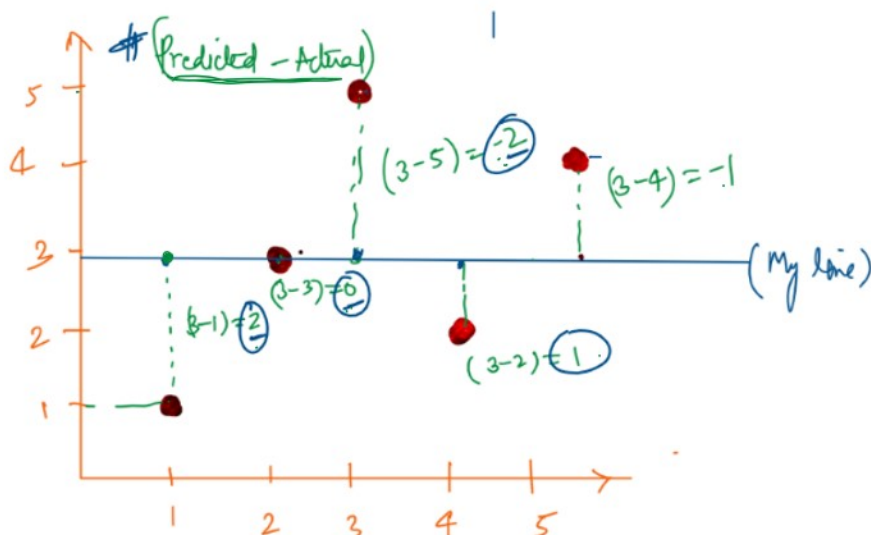
$$y = \beta_0 + \beta_1 x + u$$

Actual value Predicted value slope/gradient

$\hat{\beta}_0, \hat{\beta}_1$: derive it as well.

$$y = \underline{2} + \underline{0.5}x$$

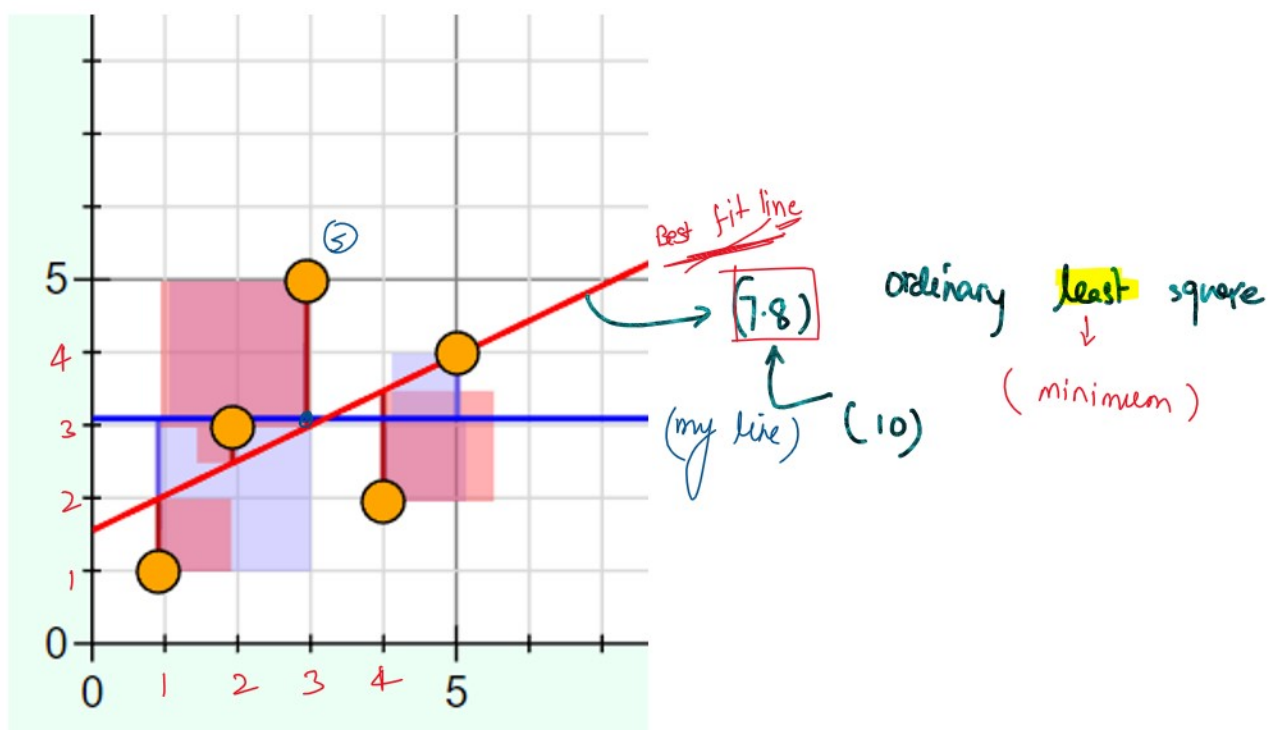




① $\sum \text{residuals} = \underline{2} + 0 + \underline{-2} + \underline{1} + \underline{-1} = 0$

(zero sum of residuals)

↓
(zero error in the model)



Observation: Residuals nullify each other and may not be truly representing the residuals.

Observation: Residuals nullify each other and may not be truly representing the residuals!

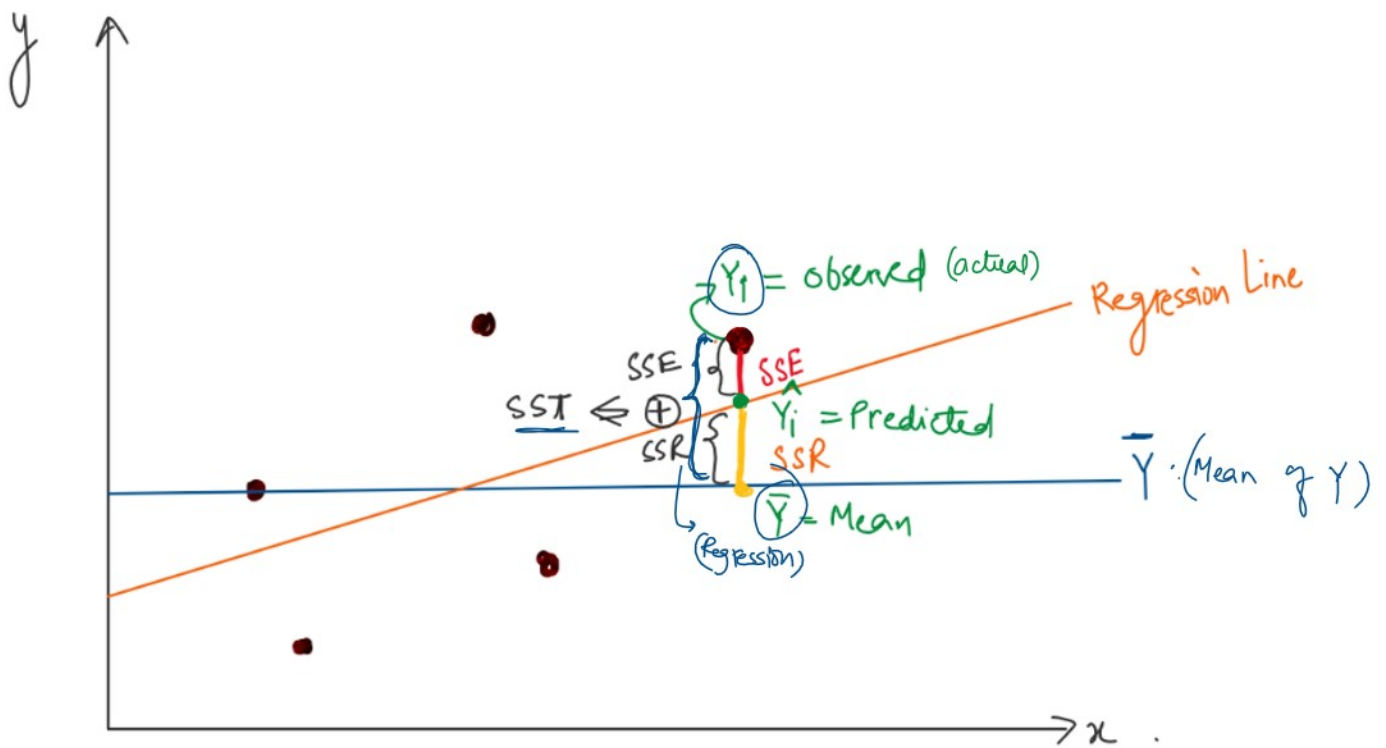
OLS Method: ordinary least squares

Sum of squared residuals

$$\text{Sum of residuals} = \cancel{2} + 0 + \cancel{-2} + \cancel{1} + \cancel{-1} = \underline{\underline{0}}$$

$$\begin{aligned}\text{Sum of squared residuals} &= 2^2 + 0^2 + (-2)^2 + 1^2 + (-1)^2 \\ &= 4 + 0 + 4 + 1 + 1 \\ &= \underline{\underline{10}}\end{aligned}$$

SSE, SSR and SST



sum of squares error (residuals) (SSE)

— it is the sum of the squared difference between the observed (actual) value and predicted value (\hat{Y}_i)

— it shows the **unexplained variance** by regression.

$$SSE = \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

sum of squares regression (SSR)

— it is the sum of the squared differences between predicted value (\hat{Y}_i) and the mean of the dependent variable (\bar{Y})

— it is a measure that describes how well

- it is a measure that describes how well our line fits the data.

$$SSR_{\text{(regression)}} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

↙ (regression)
↘ Mean

sum of squares total (SST)

- it is the squared differences between observed dependent variable and its mean.

$$SST / TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- it is a measure of total variability of the dataset.

Mathematically,

$$SST = SSR + SSE$$

$$\Rightarrow \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

↓
Total variability
of the dataset
(SST)

↓
variability explained
by the regression line
(SSR)

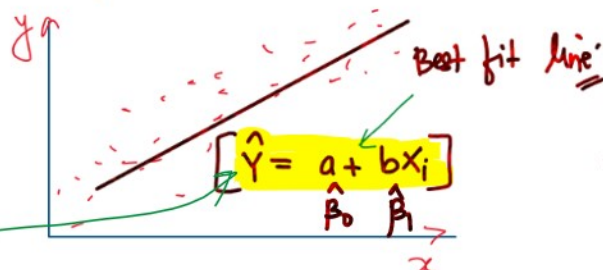
↓
unexplained
variability
(SSE)

$$100\% = 75\% + 25\%$$

↓
Accuracy

Derivation of linear regression equation

Given a set of n points (x_i, y_i) on a scatter plot.



$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$\hat{Y} = a + b x_i$$

find the best fit line; $\hat{Y} = a + b x_i$

such that the sum of squared errors in Y : $\sum_{i=1}^n (y_i - \hat{Y}_i)^2$ is minimized

(actual) (predicted)

(observed)

$$S = \sum_{i=1}^n (y_i - \hat{Y}_i)^2 = \sum_{i=1}^n [y_i - (a + b x_i)]^2$$

$$S = \sum_{i=1}^n (y_i - a - b x_i)^2$$

Using maxima/minima concept, let us partially differentiate with 'a' and 'b' respectively.

$$\left[\begin{array}{l} z = x^2 y \\ \frac{\partial z}{\partial x} = 2xy \quad \left| \quad \frac{\partial z}{\partial y} = x^2 \right. \\ \text{-(y constant)} \quad \quad \quad \text{(x constant)} \end{array} \right] \quad \boxed{\frac{\partial S}{\partial a} = 0}$$

$$\frac{\partial S}{\partial a} = \sum_{i=1}^n 2 (y_i - a - b x_i) \times (0 - 1 - 0)$$

$$\frac{\partial S}{\partial a} = -2 \sum_{i=1}^n (y_i - a - b x_i) \quad \left(\begin{array}{l} \text{(chain rule)} \\ \frac{d}{da} (y_i - a - b x_i) \\ = (0 - 1 - 0) \end{array} \right)$$

Differentiation

$$y = x^2 \quad \left\{ \begin{array}{l} \frac{dy}{dx} = 2x \end{array} \right.$$

$$y = x^3 \quad \frac{dy}{dx} = 3x^2$$

$$y = c \quad \frac{dy}{dx} = 0$$

$$y = x^n$$

$$\frac{dy}{dx} = n x^{n-1}$$

<https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/partial-derivatives/v/partial-derivatives-introduction>

$$\frac{dy}{dx} = nx^{n-1}$$

$$\frac{\partial Q}{\partial a} = -2 \left[\sum_{i=1}^n y_i - \sum_{i=1}^n a - b \sum_{i=1}^n x_i \right] = 0$$

$$\frac{\partial Q}{\partial a} = -2n\bar{y} + 2na + 2bn\bar{x} = 0$$

Dividing by $2n$

$$\Rightarrow -\bar{y} + a + b\bar{x} = 0$$

$$\Rightarrow a = \bar{y} - b\bar{x}$$

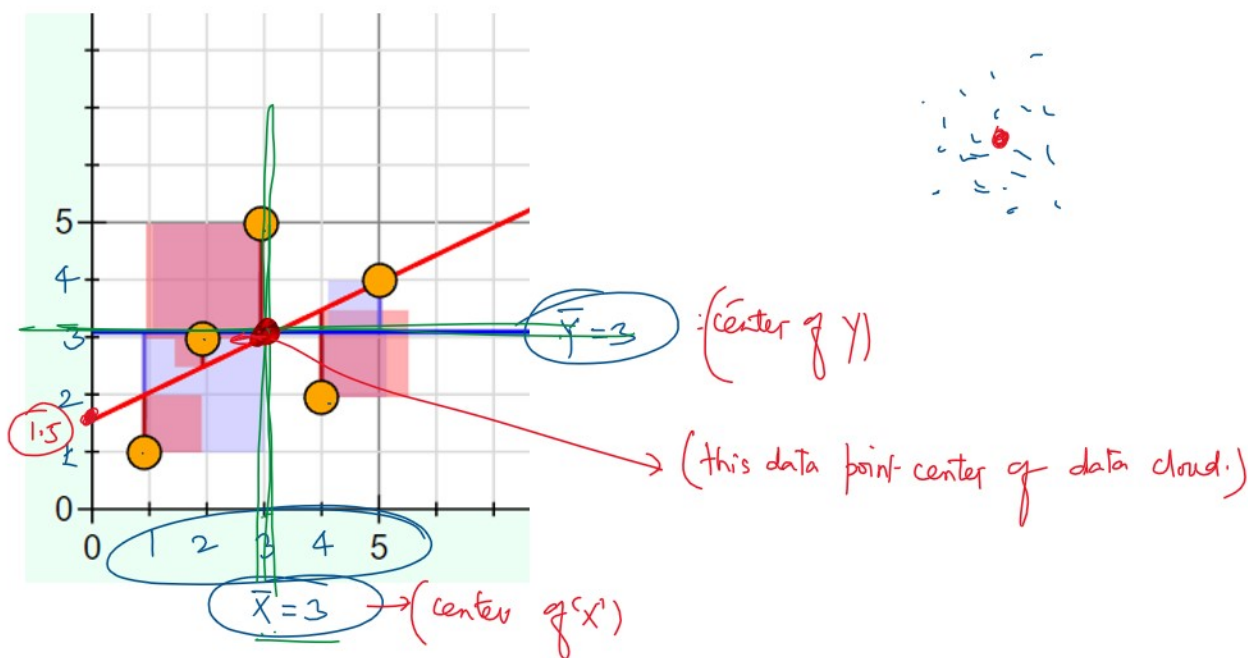
$\hat{\beta}_0$ $\hat{\beta}_1 = ??$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\Rightarrow \sum_{i=1}^n x_i = n\bar{x}$$

$$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} = a + b \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

a is constant (y -intercept) is such that the line must go through the mean of x and y



https://www.youtube.com/watch?v=4b4MUyve_U8 : Andrew Ng - DO NOT WATCH NOW -- try it after our course is over

second condition for minimizing \mathcal{Q} is:

$$\boxed{\frac{\partial \mathcal{Q}}{\partial b} = 0}$$

$$\mathcal{Q} = \sum_{i=1}^n (y_i - a - b x_i)^2$$

$$\frac{\partial \mathcal{Q}}{\partial b} = \sum_{i=1}^n 2 (y_i - a - b x_i) x_i (0 - 0 - x_i)$$

$$\frac{\partial \mathcal{Q}}{\partial b} = (-2) \sum_{i=1}^n (y_i - a - b x_i) (x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n [y_i x_i - a x_i - b x_i^2] = 0$$

let us substitute 'a' from $a = \bar{y} - b \bar{x}$

$$\Rightarrow \sum_{i=1}^n [y_i x_i - (\bar{y} - b \bar{x}) x_i - b x_i^2] = 0$$

$$\Rightarrow \sum_{i=1}^n [y_i x_i - \bar{y} x_i + b \bar{x} x_i - b x_i^2] = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i - \bar{y} x_i) + b \sum_{i=1}^n (\bar{x} x_i - x_i^2) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i - \bar{y} x_i) - b \sum_{i=1}^n (x_i^2 - \bar{x} x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i - \bar{y} x_i) = b \sum_{i=1}^n (x_i^2 - \bar{x} x_i)$$

$$\begin{array}{l} \sum_{i=1}^n b = b \sum_{i=1}^n 1 \\ \downarrow \qquad \qquad \qquad \searrow \\ b + b + \dots + n \text{ times} \quad b(1 + 1 + \dots + 1) \\ (bn) \qquad \qquad \qquad \underline{bn} \end{array}$$

$$\Rightarrow b = \frac{\sum_{i=1}^n [y_i x_i - \bar{y} x_i]}{\sum_{i=1}^n [x_i^2 - \bar{x} x_i]}$$

sum product (x_i y_i)

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\Rightarrow b = \frac{\sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i^2 - \bar{x} \sum_{i=1}^n x_i} \quad (n\bar{x})$$

$$\Rightarrow b = \frac{\sum_{i=1}^n x_i y_i - n\bar{y}\bar{x}}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}$$

Intuitively by using two expressions:

$$\sum_{i=1}^n (\bar{x}^2 - x_i \bar{x}) = 0$$

$$= n\bar{x}^2 - \sum_{i=1}^n x_i \bar{x}$$

$$= n\bar{x}^2 - \bar{x} \cdot \left(\sum_{i=1}^n x_i \right) \rightarrow n\bar{x}$$

$$= \cancel{n\bar{x}^2} - \cancel{n\bar{x}}$$

$$= 0$$

$$\bar{x} = k$$

$$\bar{x}^2 = k^2$$

$$\sum_{i=1}^n k^2 = n k^2$$

$$n\bar{x}^2$$

$$\left(\sum_{i=1}^n x_i \right) \rightarrow n\bar{x}$$

$$\text{by, } \sum_{i=1}^n (\bar{y}^2 - y_i \bar{y}) = 0$$

$$b = \frac{\sum_{i=1}^n (x_i y_i - x_i \bar{y}) + \sum_{i=1}^n (\bar{x} \bar{y} - y_i \bar{x})}{\sum_{i=1}^n (x_i^2 - x_i \bar{x}) + \sum_{i=1}^n (\bar{x}^2 - x_i \bar{x})}$$

$$\sum_{i=1}^n (x_i^2 - x_i \bar{x}) + \sum_{i=1}^n (\bar{x}^2 - x_i \bar{x})$$

$$\sum_{i=1}^n (x_i - \bar{x}) + \sum_{i=1}^n (\bar{x} - x_i) = 0$$

reverse

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$b = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

$$a = \bar{y} - b\bar{x}$$



Covariance Formula

For Population

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x}) * (y_i - \bar{y})}{N}$$

$$\rho = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

For Sample

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x}) * (y_i - \bar{y})}{(N-1)}$$

