

# Linear Regression

15 December 2023 07:52

→ model is linear → linear relationship between  $Y$  &  $X_i$ 's

Supervised Learning Algorithm

→ continuous / categorical

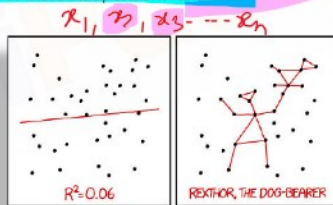


## Linear Regression using Python



### 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



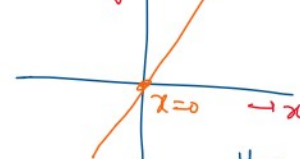
I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

www.intellipaat.com

a function in  $x$

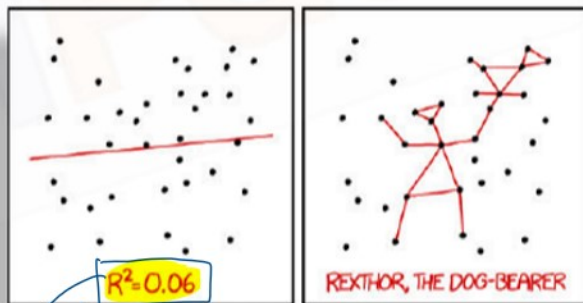
$$y = 2x$$

$$y = mx$$



$$y = 2x$$

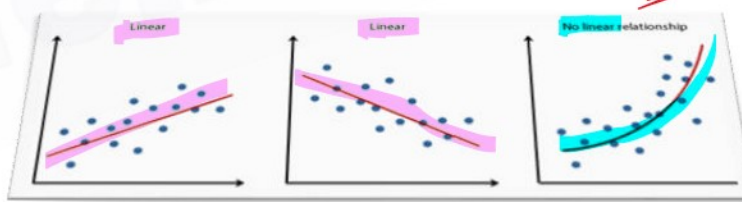
y	x
0	0
-2	-1
2	1



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



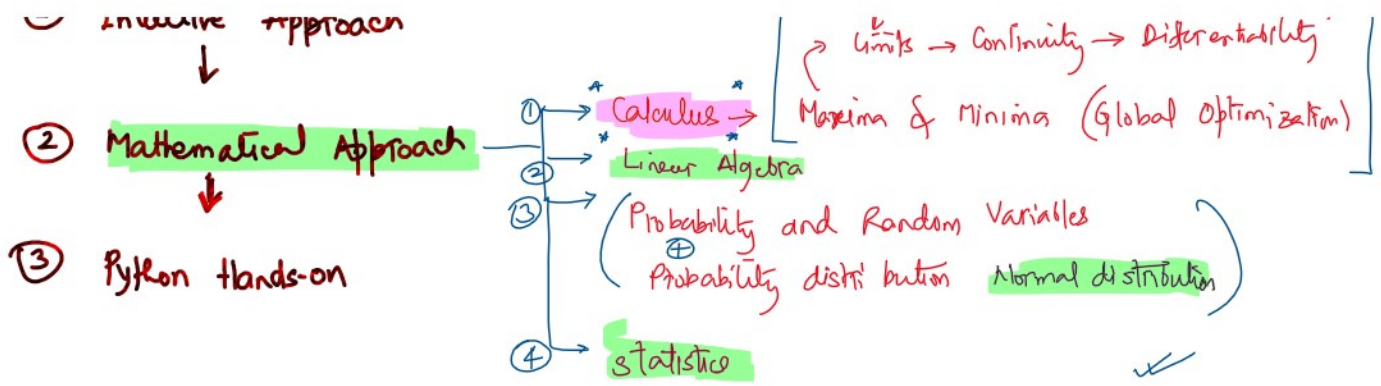
① Intuitive Approach

→ Calculus

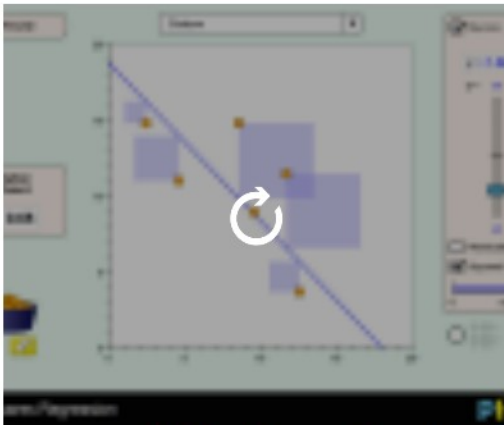
function

→ limits → Continuity → Differentiability

→ ... (f(x) is a solution to ...)



# Least-Squares Regression



mathematical equation

$$Y = 0.83X + 0.75$$

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

slope      intercept

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

**Best-Fit Line**

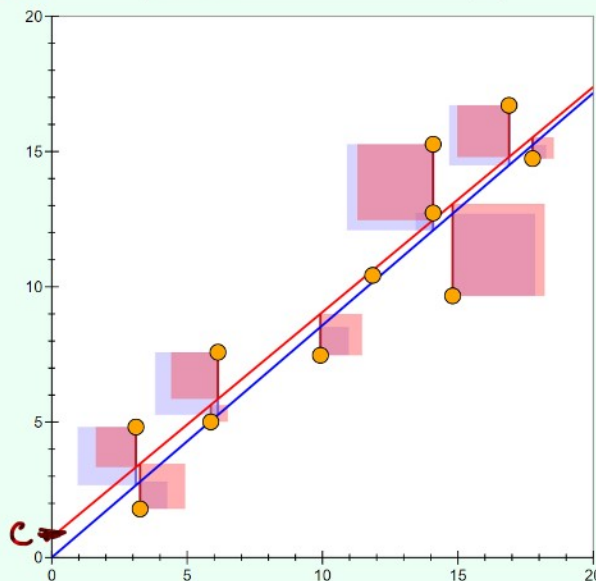
☒ Best-Fit Line

☒ Residuals

☒ Squared Residuals

0      sum

**Correlation Coefficient**



☒ My Line

$y = 0.86x + 0.00$

$y = a x + b$

☒ Residuals

☒ Squared Residuals

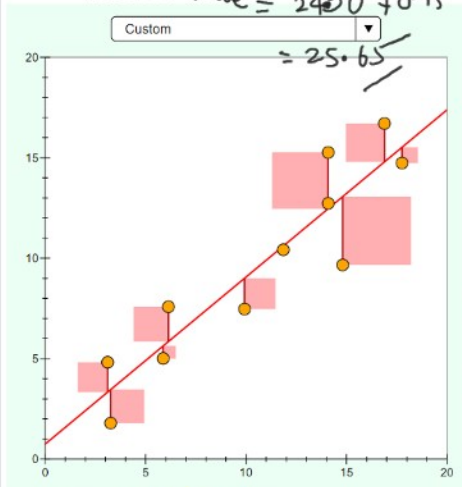
0      sum



$$y = 0.83x + 0.75$$

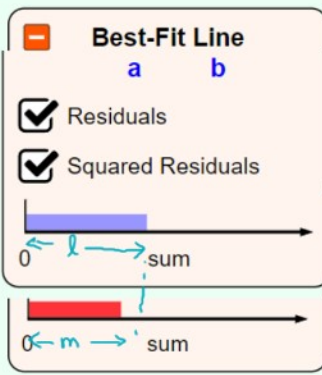
$$y_{\text{Predicted value}} = 0.83 \times 20 + 0.75$$

at  $x = 20$



• : data points

— : Linear Regression Line



l: intuitive best fit line

m: LR model best fit line

a)  $l > m$

b)  $l < m$

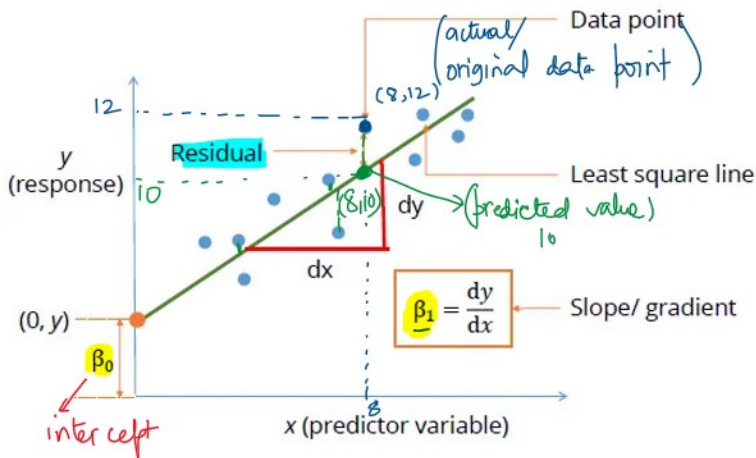
c)  $l = m$

d) Not sure

future

• data point

(Linear Regression/ Best fit/ least square line)



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

Residual

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

Actual value

Predicted value

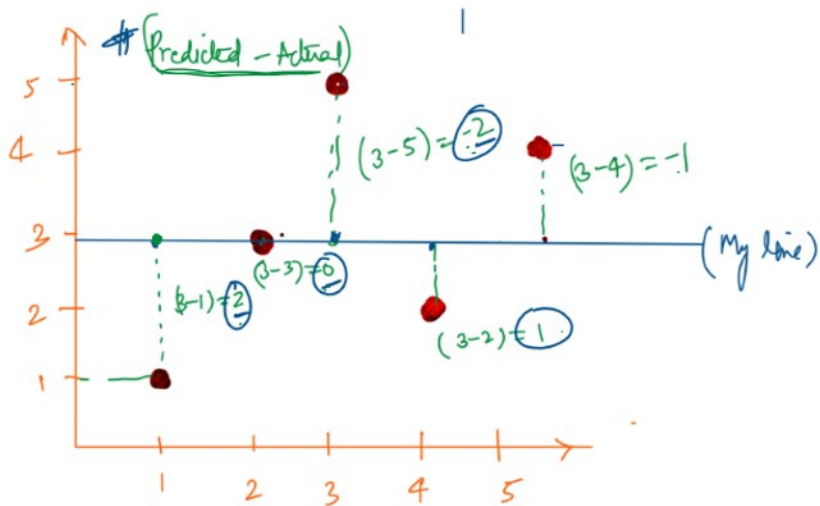
slope/ gradient

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

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

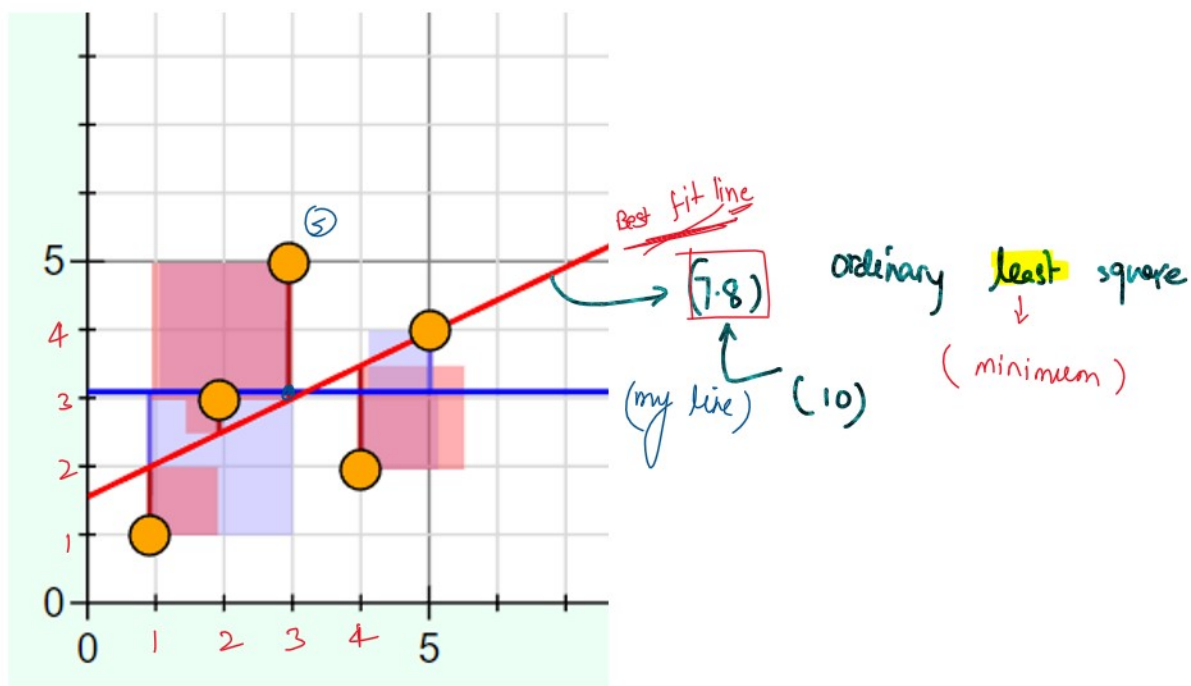


+ve residues cancel -ve residues



①  $\text{sum of residues} = \cancel{2} + 0 \cancel{-2} + 1 \cancel{-1} = 0$

(zero sum of residues)  
 $\downarrow$   
 (zero error in the model)



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

OLS Method: ordinary least square

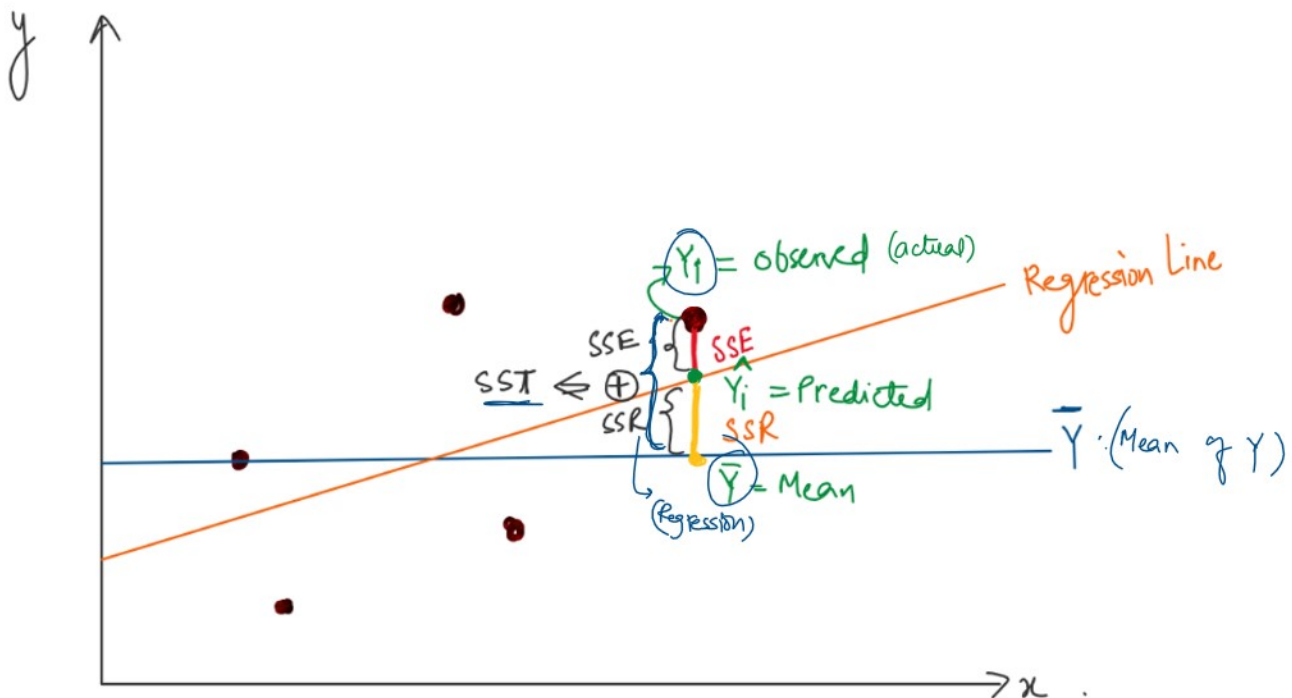
## OLS Method: ordinary least squares

sum of squared residues

$$\text{sum of residues} = \cancel{-2} + 0 + \cancel{-2} + \cancel{+1} + \cancel{+1} = \underline{\underline{0}}$$

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

## # SSE, SSR and SST



## # sum of squares error (residues) (SSE)

- it is the sum of the squared differences between the observed (actual) value and



- it is the sum of the squared differences 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 our line fits the data.

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

(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,

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)

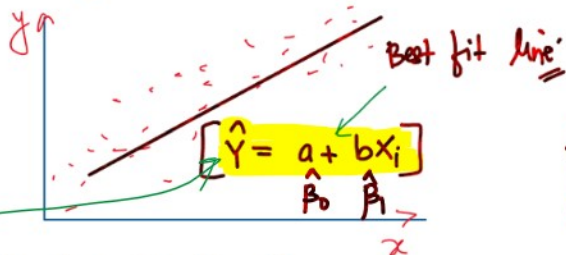
unexplained  
variability  
(SSE)

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

↓  
Accuracy

# Derivation of linear regression equation

Given a set of 'n' points  $(X_i, Y_i)$  on a scatter plot.



$$\hat{Y} = \beta_0 + \beta_1 X$$

$$\hat{Y} = a + bX_i$$

Find the best fit line;  $\hat{Y} = a + bX_i$

such that the sum of squared errors in  $Y$ :

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

actual (Predicted)  
(Observed)

$$Q = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n [Y_i - (a + bX_i)]^2$$

$$Q = \sum_{i=1}^n (Y_i - a - bX_i)^2$$

Using maxima/minima concept, let us partially differentiate

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

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

$$\frac{\partial Q}{\partial a} = \sum_{i=1}^n 2 \underbrace{(y_i - a - bx_i)}_{(y \text{ constant})} \times (0 - 1 - 0)$$

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

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

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

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

(na)

Dividing by 2n

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

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

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

Differentiation

$$\begin{array}{l} y = x^2 \\ \frac{dy}{dx} = 2x \end{array} \quad \left\{ \begin{array}{l} y = x^3 \\ \frac{dy}{dx} = 3x^2 \end{array} \right. \quad y = c \quad \frac{dy}{dx} = 0$$

$$y = x^n$$

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

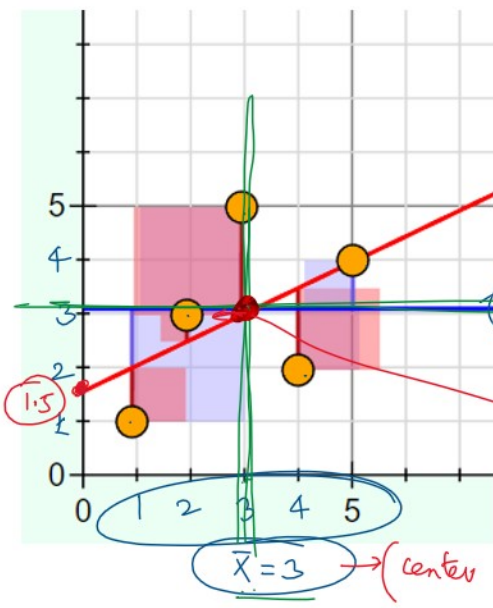
$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

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

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

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





(center of Y)

(this data point-center of data cloud.)

(center of X)

$$a = 1.5$$

[https://www.youtube.com/watch?v=4b4MUyve\\_U8](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  $\mathcal{J}$  is:

$$\frac{\partial \mathcal{J}}{\partial b} = 0$$

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

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

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

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

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

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

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

$$\Rightarrow \sum_{i=1}^n \left[ \underbrace{y_i x_i - \bar{y} x_i}_{\text{red}} + \underbrace{b \bar{x} x_i - b x_i^2}_{\text{red}} \right] = 0$$

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

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

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

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

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

sum product (x<sub>i</sub> y<sub>i</sub>)

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

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

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

Intuitively by using two expressions:

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

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

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

$$= n\bar{x}^2 - n\bar{x}^2 = 0$$

$$\bar{x} = k$$

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

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

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

$$\begin{aligned}
 &= n\bar{x}^2 - \bar{x} \cdot \left( \sum_{i=1}^n x_i \right) \rightarrow (n\bar{x}) \\
 &= n\bar{x}^2 - n\bar{x}^2 \\
 &= 0
 \end{aligned}$$

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

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

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

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

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

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



## 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

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

## # $R^2$ : Coefficient of determination

\* It is a statistical measure of how well the regression line approximates the actual data.

\* It is a measure about the goodness of fit of a model.

$$SST = SSR + SSE$$

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

Total variability  
of the dataset  
(SST)

variability explained  
by the regression line  
(SSR)

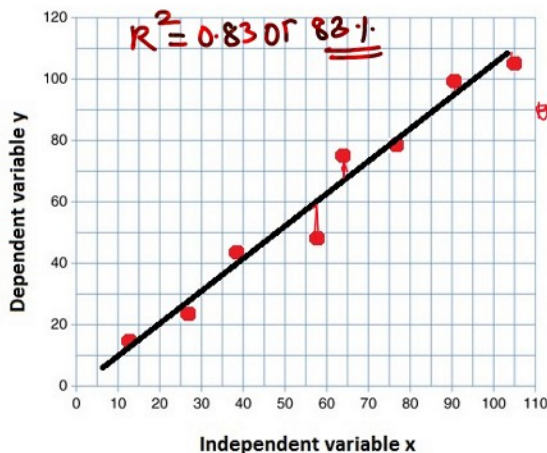
unexplained  
variability  
(SSE)

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

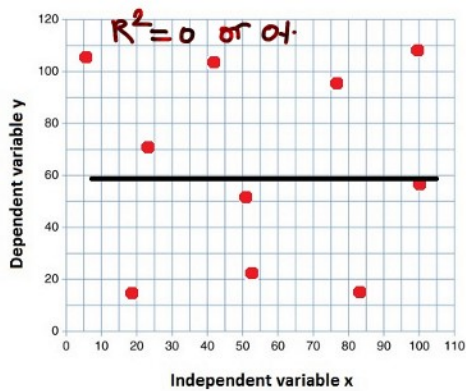
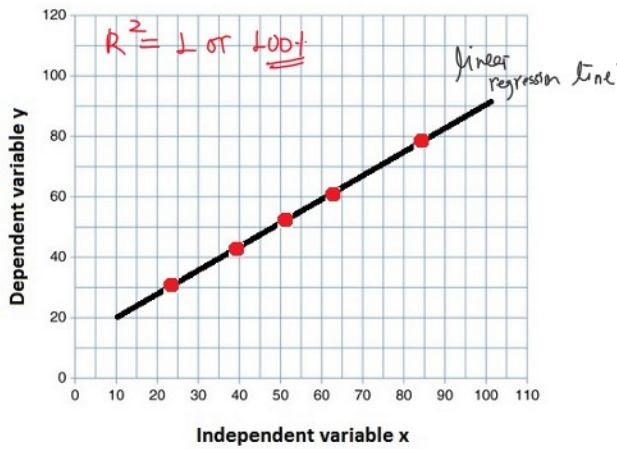
$$R^2 = 1 - \frac{\text{sum squared error (SSE)}}{\text{total sum of squares (SST)}}$$

$$R^2 = \frac{\text{sum of squared regression (SSR)}}{\text{total sum of squares (SST)}}$$

Let us say  $R^2 = 83\%$  which means that 83% of the variation in  $y$  values is accounted for by the  $x$  values.



17% of total 'y' variance is unexplained variability.



R-squared or coefficient of determination is a measure that gives proportion of variation in target variable (y) explained by the linear regression model.

$$\uparrow\uparrow R^2 = \left( 1 - \frac{SSE \downarrow\downarrow}{SST} \right)$$

Given regression line is getting very close to actual data points, sum of squares error (SSE) decrease and hence  $R^2$  increases

$$\downarrow\downarrow R^2 = \left( 1 - \frac{SSE \uparrow\uparrow}{SST} \right)$$

### # Problem with $R^2$ statistic

$R^2$  value never decreases no matter the no. of variables we add to our regression

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots$$



$R^2$  value never decreases no matter the no. of variables we add to our regression model.

$R^2$  either remains the same or increases with the addition of new independent variables.

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

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 x_3 + \dots + \hat{\beta}_n x_n$$

↓  
75%

— MLR

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 \quad \text{--- } 0.75 \text{ or } \underline{75\%} \quad \overset{R^2}{\downarrow}$$

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \quad \text{--- } 85\% \quad \rightarrow \text{'Inflation in accuracy'}$$

## # Adjusted $R^2$ :

Unlike the std.  $R^2$ , which simply tells you the proportion of variance explained by the model,

Adjusted  $R^2$  takes into account the no. of predictors ( $x_i$ 's), independent variables in the model

$$\text{Adjusted } R^2 = \left\{ 1 - \frac{[(1 - R^2)(n - 1)]}{(n - K - 1)} \right\}$$

where  $n$ : represents the no. of data points

$K$ : " " no. of variables

$R^2$ : std.  $R^2$  value.

$$(101 - 1) = 100$$

$$\left[ 1 - \frac{(1 - R^2)(n - 1)}{(n - K - 1)} \right]$$

↓  
denominator

$$R^2 = 0.85$$

$$(1 - R^2) = 0.15$$

$$(1 - R^2)(n - 1) = 0.15 \times 100 = 15$$

$$(n - K - 1) = (101 - 1 - 1) = 99$$

$$\text{Adj. } R^2 = \left( 1 - \frac{15}{99} \right)$$

$$1 - (15/99) = 0.8485 \approx 85\%$$

$$K = 10$$

$$\left( \frac{15}{101 - 10 - 1} \right) = \left( \frac{15}{90} \right)$$

$$\left( \frac{15}{101-10-1} \right) = \left( \frac{15}{90} \right)$$

$$\text{Adj } R^2 = \left( 1 - \frac{15}{90} \right)$$

$$1 - (15/90) = 0.8333$$

83.33 ↓↓

over

"random independent variable" 97.1%

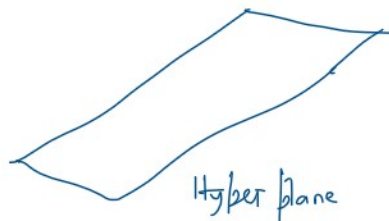
OLS Regression Results

Dep. Variable:	marks	R-squared:	0.971
Model:	OLS	Adj. R-squared:	0.967
Method:	Least Squares	F-statistic:	265.0
Date:	Sat, 30 May 2020	Prob (F-statistic):	2.04e-07
Time:	23:41:47	Log-Likelihood:	-17.372
No. Observations:	10	AIC:	38.74
Df Residuals:	8	BIC:	39.35
Df Model:	2		
Covariance Type:	nonrobust		

OLS Regression Results

Dep. Variable:	marks	R-squared:	0.971
Model:	OLS	Adj. R-squared:	0.962
Method:	Least Squares	F-statistic:	116.1
Date:	Sat, 30 May 2020	Prob (F-statistic):	4.28e-06
Time:	23:44:34	Log-Likelihood:	-17.364
No. Observations:	10	AIC:	40.73
Df Residuals:	7	BIC:	41.64
Df Model:	2		
Covariance Type:	nonrobust		

$$\hat{Y} = 0.25 + 1.5X_1 + 2X_2$$



overfitting vs underfitting  
bias vs variance

## # Multicollinearity

Causes, effects and detection using VIF

↳ (variance inflation factor)

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$$

obese weight  
BMI  
Height

$$BMI = f(H, W)$$

$$X_2 = f(X_1, X_3)$$