

MATHS - I : Statistics

- Statistics is the branch of mathematics dealing with the Data Collection, Data Analysis, Interpretation, Data presentation, organizing the Numerical Data.

Statistics = measurement + Analysis.

Statistics

Inferential Statistics

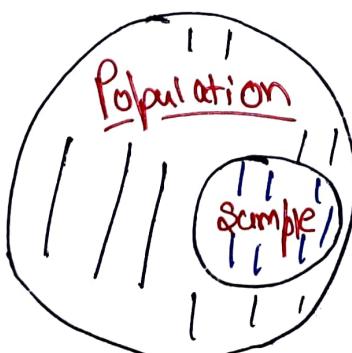
- Many times, a collection of entire data is impossible, Hence a subset of data is collected
- from the subset we get conclusions about the entire population. This is Inferential stats.

Descriptive Statistics

- These are used to summarize data, such as the mean, standard deviation & etc.
- While applying statistics on data we can find underlying relationships in between the variables.

* Population:

- Totality of all values or complete list of observations.
- All data points about the subject under study.
→ Can be people, vehicles, sales, cats, houses etc.



Taking sample from population using Pandas (Python)

→ import Panda as pd.

```
df = pd.read_csv('Sales1.csv')
```

```
print(df.sample())
```

sample() is a method.
and it always give
Random values.

* Sample: A sample is a subset of a population, usually a small portion of the population that is being analysed. It is expensive to perform an analysis on an entire population. So, by analysing sample helps to draw the conclusions about a population.

* Variable: In statistics what we are examining is called variable which is categorized, countable and measurable.

Eg: people have different heights, weights and professions.

→ Variable represents the characteristic representation of study.

→ Variable is an attribute, has a value, and can store any data type.

→ Types of variable

Quantitative DATA

- Data that is measured in Numbers. It deals with numbers that make sense to perform arithmetic calculations.
- Eg: Height, weight, etc.

① Quantitative variables:

- Quantitative variables are numeric
- we can perform arithmetic calculations.
- Represent a measurable quantity.
- Eg: Height, weight etc.

② Categorical or Qualitative variables:

- Stores values which are names or labels.
- Eg: Type of person: male or female.
Review about food: good, bad, okay.
color of bike: red, green, blue.

Types of Categorical variable:

Ordinal

- logical order is possible to do analysis.
- Eg: Grades in exams

Nominal

- There is no logical ordering with respect to actual values.
- Eg: Hair colour.

Categorical DATA

Refers to the value that place "things" into different groups or categories.

- Eg: Hair colour, type, letter grade.

Measures of SPREAD

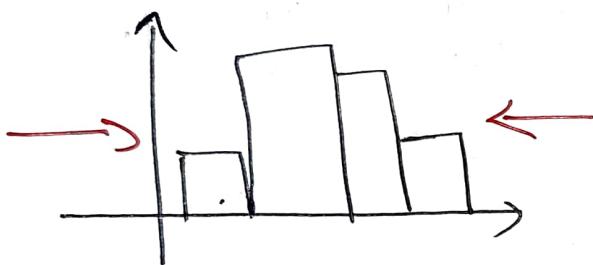
⇒ Types of Quantitative variables:

Discrete Variables

- Whole Numbers
- No. of Pets Owned
- No. of people in family
- No. of bikes or cars.

Continuous Variables

- float numbers
- weight
- salary
- Bank balance.



measures of Centre

- mode
- median
- mean

measures of spread

- Range
- Standard Deviation.

① mode → most frequently observed value.

Eg: 154, 139, 154, 192, 180, 140, 154, 155, 192

mode → 154.

② median → value in the middle of an ordered dataset.

$$\Rightarrow \boxed{\frac{n+1}{2}}$$

③ mean → Average

$$\text{mean} = \frac{\sum x_i}{n}$$

summation of all values
Total number of values.

$$\text{Sample mean} = \bar{x} = \frac{\sum x_i}{n}$$

Measures of SPREAD:

* RANGE:

Range means difference in between minimum values and maximum value. It explains about the data is in between min and max values.

* Standard Deviation:

The Standard Deviation is a measure of how spread out numbers.

→ Square root of variance.

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

x_i = observations

\bar{x} = sample mean

n = total no. of observations

STD tells us how close the values in a Data set are to the mean.

* Variance:

- Average of squared differences from the mean.
- Variance is the average of squared differences from the mean.
- How far the data points are in a population from the population mean.

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

Outliers

A value which is very far from common Data.
→ can be largest value or the smallest.

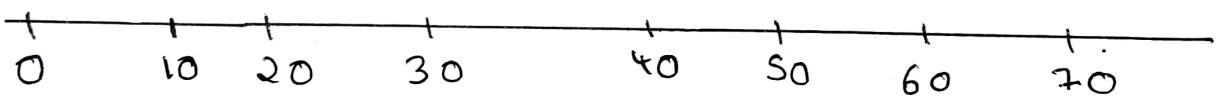
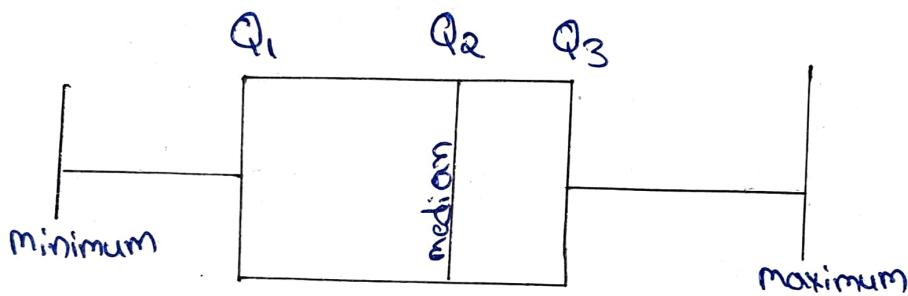


The Five Number Set

1, 2, 3, 4, --- Q1 --- 50 --- 98, 99, 100.



Box - PLOT:



* IQR: Inter quartile Range:

$$\text{IQR} = Q_3 - Q_1.$$

- ① Box plot shows the maxi, mini, medium, first Q and third Quartile of the DataSet.
- ② gives the overall & statistical information about data distribution.
- ③ Good for detecting the Outliers.

Eg:

```
import matplotlib.pyplot as plt  
data = [10, 20, 99, 30, 140, 50]  
plt.boxplot(data)  
plt.show()
```

* A Data Value is considered to be an outlier if:

Data Value $< Q_1 - 1.5(\text{IQR})$

Small Outlier

OR

Data Value $> Q_3 + 1.5(\text{IQR})$

Large Outlier.

Eg:

Q_1	Q_2	Q_3
10	25	33

10, 11, 12, 25, 25, 27, 31, 33, 34, 34, 35, 36, 50, 59

$$\text{IQR} = Q_3 - Q_1 = 36 - 25 = 11$$

Lower Whisker = $Q_1 - 1.5(\text{IQR})$
= $25 - 1.5(11)$
= 8.5

Upper Whisker = $Q_3 + 1.5(\text{IQR})$
= $36 + 1.5(11)$
= 52.5

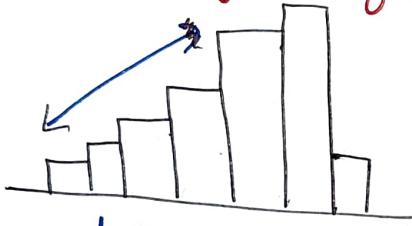
∴ Any value less than 8.5 or greater than 52.5 is an outlier

Symmetry and Skewness:



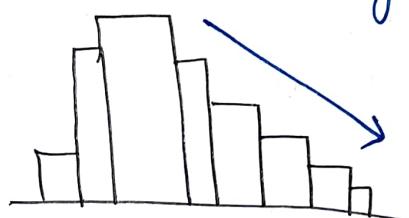
Symmetry

Same shape on both sides of the median.



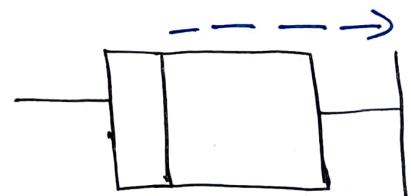
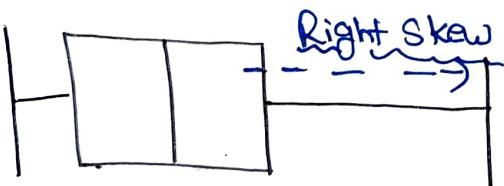
Left skewed

Long tail that trails towards left



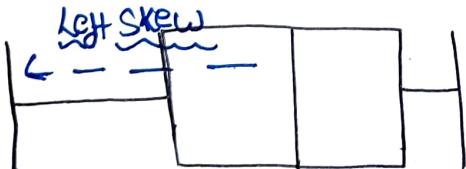
Right skewed

long tail that trails toward Right



→ If Right side of the box is large so Right Skew box plot.

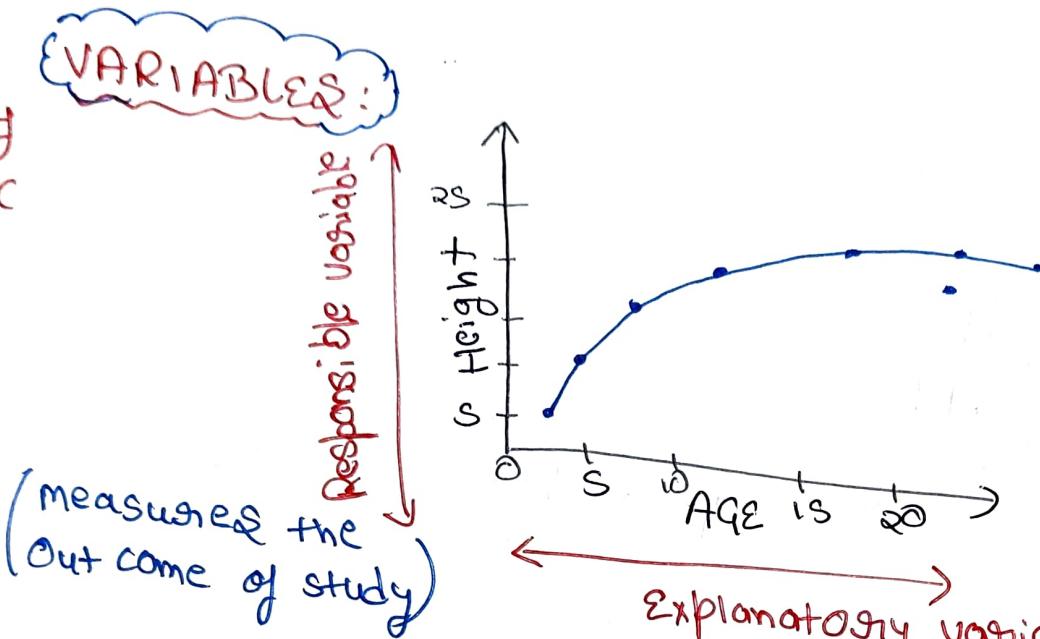
→ If box is equal then, go for the large whisker.



Symmetric

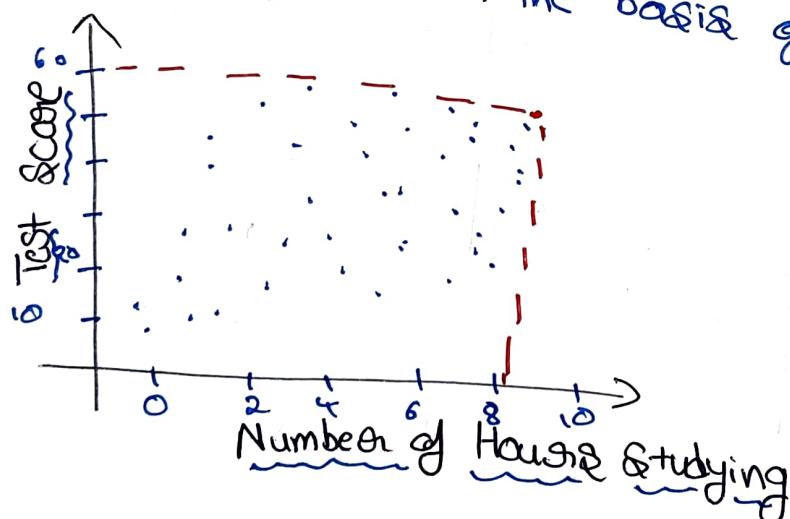
(i) IV
C

VARIABLES:



Scatter Plot:

Eg: State of growth / loss on the basis of one variable.



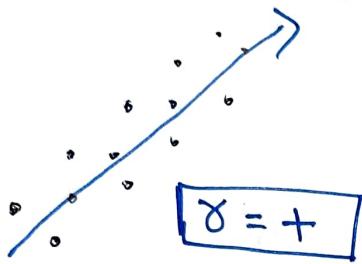
⇒ If two variables have any relation then we can find this relationship by using scatter plot.

Correlation:

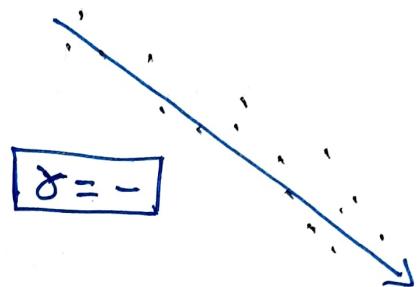
Explains about Direction and Strength of Relationship between two variables.

(6)

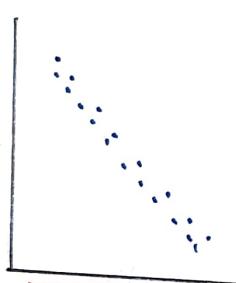
(i) If upward than
Correlation is positive



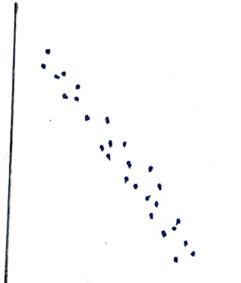
(ii) Downwards than
Correlation is negative.



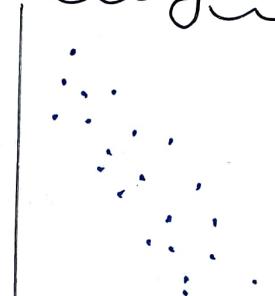
Direction & strength



$r = -1$
Perfectly
Negative



$r = -0.7$

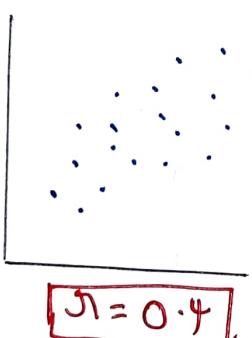


$r = -0.4$

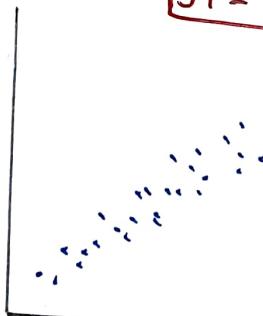


No relation

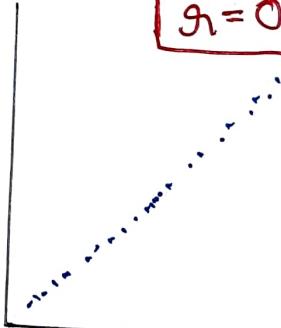
$r = 0$



$r = 0.4$



$r = 0.7$



$r = 1$

Perfectly Linear
positive Relation.

$$r = \frac{1}{(n-1)s_x s_y} \sum (x_i - \bar{x})(y_i - \bar{y})$$

understanding Correlation:

$$r = \frac{1}{(n-1)s_{xy}} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$n \rightarrow$ Number

of observations.

$s \rightarrow$ Standard Deviation (for x and y)

Eg:

A teacher wants to determine the Corr between the number of hours spent studying and test scores.

2

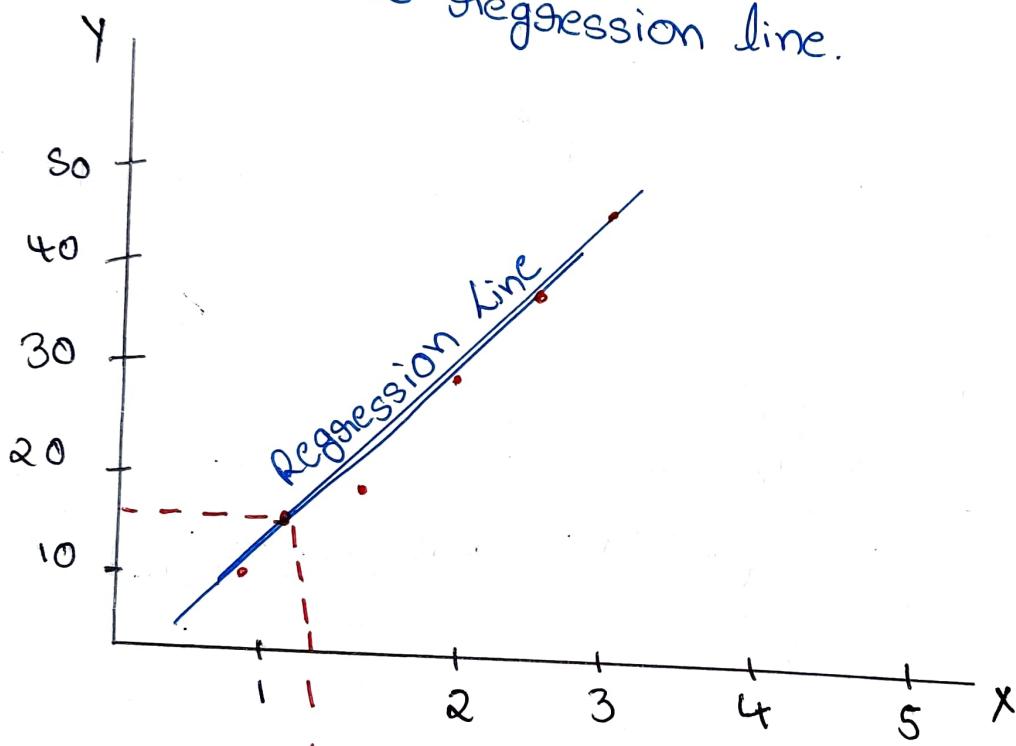
Student	Time Spent x_i	Score (100) y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$
John	13	53	$13 - 12.5$ = 0.5	$53 - 68$ = -15	$(0.5)(-15)$ = -7.5
Allie	15	69	$15 - 12.5$ = 2.5	1	2.5
Mark	7	92	-5.5	24	-13.5
Samantha	3	10	-9.5	-8	82.5
Jessica	10	89	-2.5	17	42.5
Joseph	27	99	14.5	31	449.5
	$\bar{x} = 12.5$	$\bar{y} = 68$			<u>Sum = 821</u>
	$s_x = 8.28$	$s_y = 32.91$			

$$r = \frac{1}{(6-1)(8.28)(32.91)} (821) \Rightarrow r = 0.602$$

\Rightarrow The correlation is +0.6, its upward direction.

Regression

- Regression explains about how we can draw line in between the points
- The line represents the pattern of the data.
- Line is called regression line.



we can predict any point using this line when both variables have Relation.

- Regression line predicts the change in Y when X increases by one unit.
- Here change in Y can be increase / decrease

Regression line:

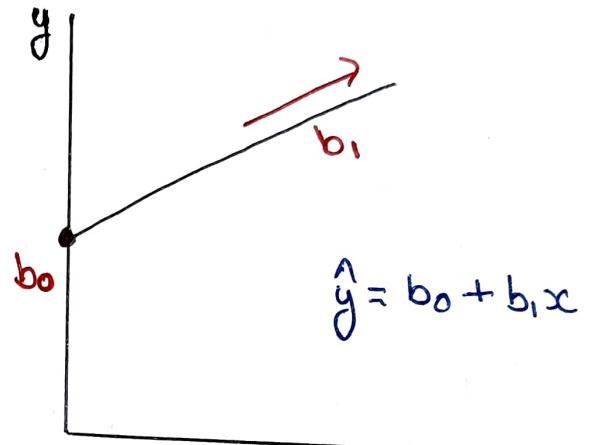
$$\hat{y} = b_0 + b_1 x$$

↓ ↓ ←
 Predicted value of Y Y intercept Any value of x.
 Slope

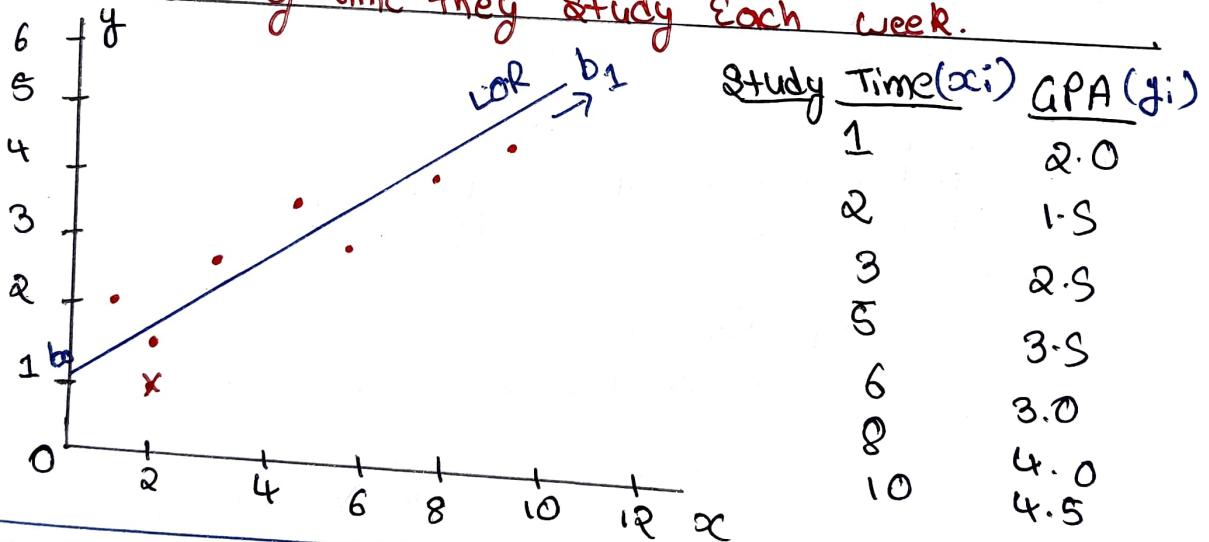
$$\hat{y} = b_0 + b_1 x$$

$$b_1 = g_1 \frac{\sum y}{\sum x}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$



Q. Suppose A Researcher wants to Predict A Student's GPA from the Amount of time they Study Each week.



$$\therefore \bar{x} = 5, \bar{y} = 3, \sum x = 32.6, \sum y = 30.8, g_1 = 0.94$$

we know,

$$\hat{y} = b_0 + b_1 x$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_1 = g_1 \frac{\sum y}{\sum x}$$

$$b_0 = 1.45$$

$$b_1 = 0.311$$

$$\therefore \hat{y} = 1.45 + 0.311x$$

So, As study time increases by one unit, we predict a student's GPA to increase by 0.311

R (r)

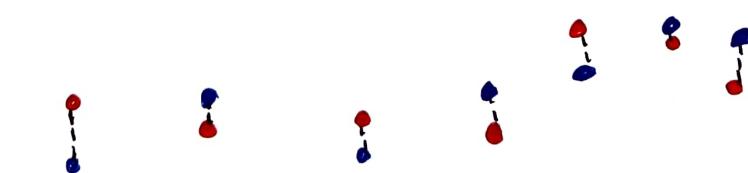
- Has values between -1 & 1
- measures the Linear relation -ship between two quantitative variables with respect and to direction and strength.

R^2 (r^2)

- Has values between 0 & 1
- Is a measure of how close each data point fits to the Regression line.
- Tells us How well the regression line predicts Actual values.

Eg:

y



- Predicted value
- Actual value.

$$r^2 = 0.90$$

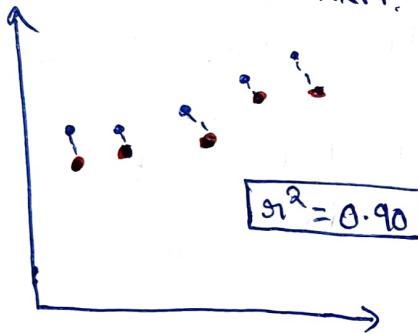
x

Case - 1

R squared value is high

Explains the actual values and predicted values are close together.

→ Actual & predicted values having very less distance b/w them.

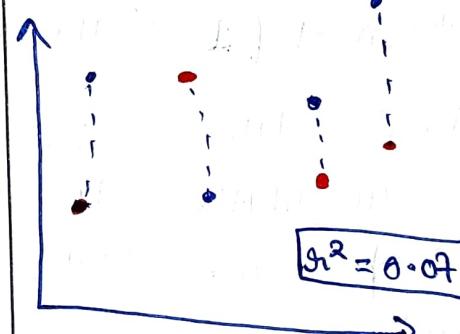


Case - 2

r^2 squared value is low

Explains the regression line which the data points are does not fit well.

→ values having large distance b/w them.



Case - 3

r^2 squared = 1

Perfect prediction

→ means we can predict value

of y for any given value of x.



NOTE:

When prediction (score value, threshold value, Score) $> 90+$, than its a good model / prediction.

& Score < 55 or so is a bad model.

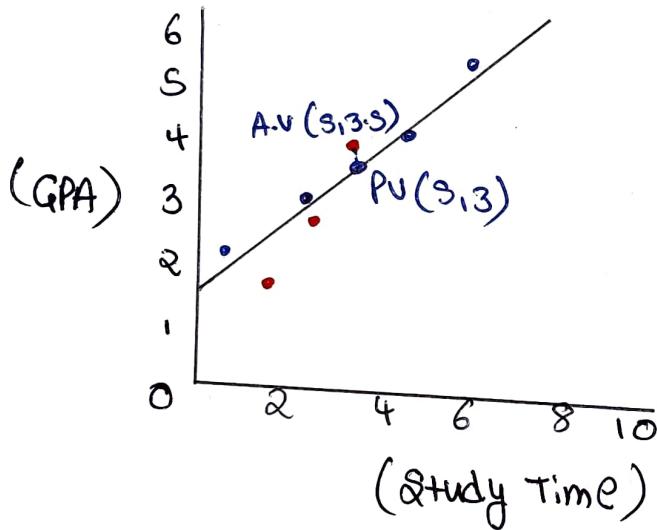
⇒ Residual:

How far the Predicted value is from the actual value, tells us the error in prediction.

Residual = Actual value of y - Predicted value of y

$$R = y_i - \hat{y}$$

Eg:



$$\begin{aligned} \text{Residual} &= 3.5 - \hat{y} \\ \hat{y} &= 1.4S + 0.311x \\ &= 2.072 \end{aligned}$$

① Residual = $3.5 - 2.072$
= 1.3 (Approx) (positive)

② Residual
for Another
value
= $1.3 - 2.072$
= -0.72
(negative).

→ Residual positive → Actual value is above of prediction is below.
→ Residual negative → prediction is high if AV is low.

NOTE → In maths it's called Residual but in ML it is called COST function (Loss / Error).

MATRIX:

Scalar

24

Vector

row $[2 \ -8 \ 7]$

column

$$\begin{bmatrix} 2 \\ -8 \\ 7 \end{bmatrix}$$

Matrix

$$\begin{bmatrix} 6 & 4 & 24 \\ 1 & -9 & 8 \end{bmatrix}$$

rows x columns.