

Statistic :-

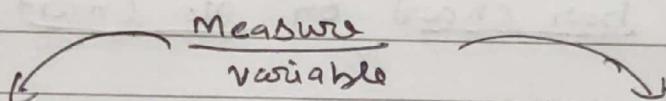
1. Inferential
2. Descriptive

1. Inferential Statistic :-

It deal with taking sample an analysis sample making judgments or claims about population

2. Descriptive Statistic :-

It refer to getting data and taking about it

Quantitative data

Data that can measure in numbers.
It deal with the number that make sense to perform Arithmetic calculations with

Quantitative variable

height
weight
midterm score

Categorical data

Refer to the value that place "thing" into different group or category

Categorical variable

Hair color
Type of cat
Letter grade

Categorical variableCategorical and ordinal

Logical ordering to the value of a categorical variable

Ex :- Letter grade

A+, A, B+, B, C+, C

Categorical and Nominal

No logical ordering to the value of a categorical variable

Ex :- Hair color

Black, Blonde, Brown

Quantitative Variable

Discrete

Refer to the variable that can only be measured in count numbers

e.g. No. of Pets

0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20

continuous

Refer to variable that can take on any numerical value

e.g. Weight

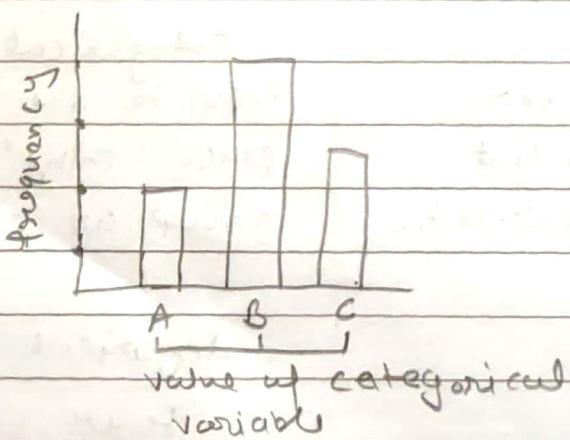
105 185 170 683

C

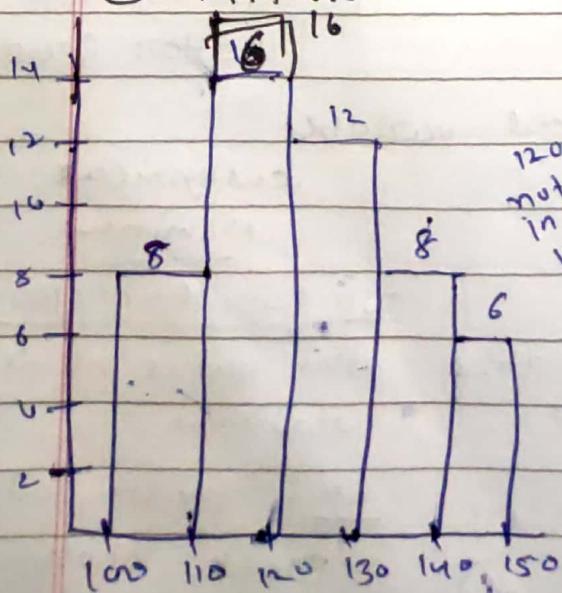
Charts

Categorical Data :-

for categorical data we use bar chart or pie chart



Quantitative data :-



Stemplot, Histogram, TimePlot

120 is included weight
not included in this interval 100 - 110

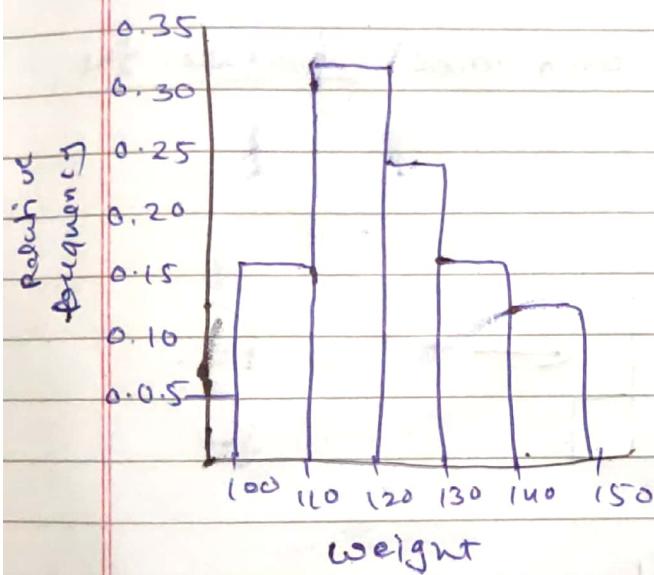
frequency
8
16
12
8
6

frequency distribution

Relative frequency distribution

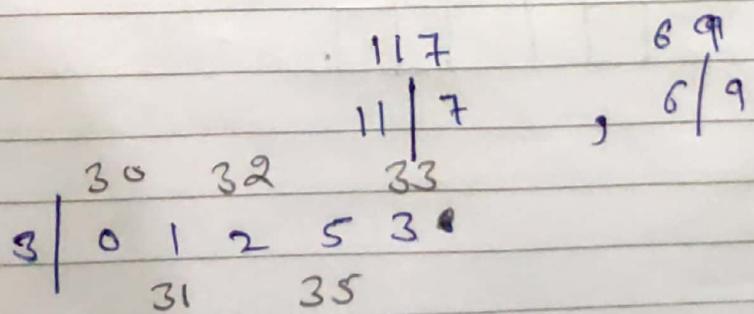
Weight	Relative frequency		Relative frequency
100 - 110	8	$8 \div 50$	0.16
110 - 120	16	$16 \div 50$	0.32
120 - 130	12	$12 \div 50$	0.24
130 - 140	8	$8 \div 50$	0.16
140 - 150	6	$6 \div 50$	0.12
	<u>50</u>		<u>1</u>

Relative frequency Histogram



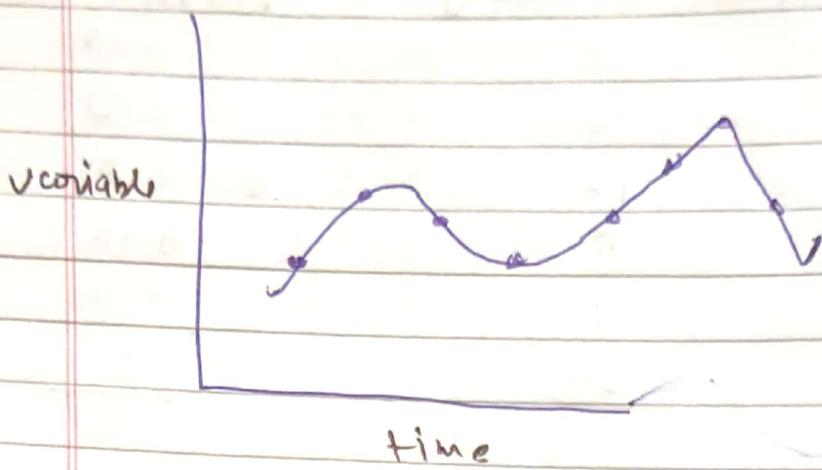
STEMPLOTS

STEMP	LEAF
Refer to the all of the other number except the last number	Refer to the very last number



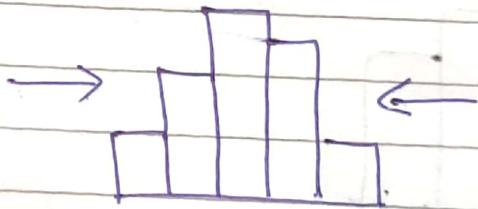
TIME PLOT

Show how a variable changes over time



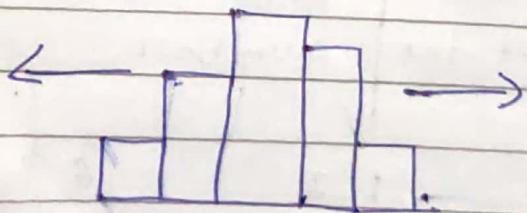
Note :-

mean, median and mode, measure the central tendency



Note :-

Range and SD measure the spread of data



measure of center :-

sample Data = 139, 140, 154, 154, 154, 154, 155, 180, 192, 192

$$\text{mode} = 154, \text{ median} = 154$$

$$\text{mean} = 165.6$$

measure of spread :-

$$\boxed{\text{Range} = \text{max} - \text{min}}$$

Sample S.D

$$\text{Range} = 192 - 139 = 53$$

$$\boxed{S.D = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}}$$

x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
10	-5.4	29.16
12	3.4	11.56
16	0.6	0.36
19	3.6	12.96
20	4.6	21.16
		<u>75.6</u>

$$S.D = \sqrt{\frac{75.6}{4}}$$

$$= \sqrt{18.8}$$

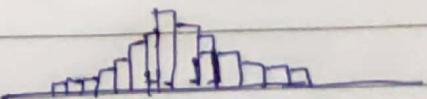
$$S.D = 4.336$$

$$\bar{x} = 15.4$$

Note :- S.D tell us how close the value in a Data Set are to the mean



Note :- Small S.D
Small amount of variability of Dataset



Note :- High S.D
High amount of variability of data set

Variance

Sample variance

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

What is the mean weight and weight S.D?

Weight of student \rightarrow 105 156 145 172 160

$$\bar{x} = 135.6 \quad S.D = 31.75$$

Suppose that during the winter, each person wear an extra 5 pounds of clothes. Under these conditions, what is the new mean and S.D?

$$\bar{x}_{\text{new}} = 135.6 + 5 = 140.6$$

S.D_{new} = 31.75 \rightarrow bcz spread of data is same for old and new data

Transforming Data

measure of center
(mean, median, mode)

Note :-
Affected by

$\oplus \ominus \times \%$

measure of spread
(range, S.D)

Note :-
Affected by

$\oplus \ominus \%$

Suppose that in order to stay hydrated, these student drink 2.5 mL of water for every pound of they weight plus 750 mL of water every day. What is the mean and S.D. for the amount of water someone every day?

$$\text{Weight} = 105 \quad 156 \quad 145 \quad 172 \quad 100$$

$$\bar{x} = 135.6$$

$$S.D = 31.75$$

$$\bar{x}_{\text{new}} = (135.6) \times 2.5 + 750$$

$$\boxed{\bar{x}_{\text{new}} = 1089}$$

$$S.D_{\text{new}} = (31.75) + 2.5$$

$$\boxed{S.D_{\text{new}} = 79.38}$$

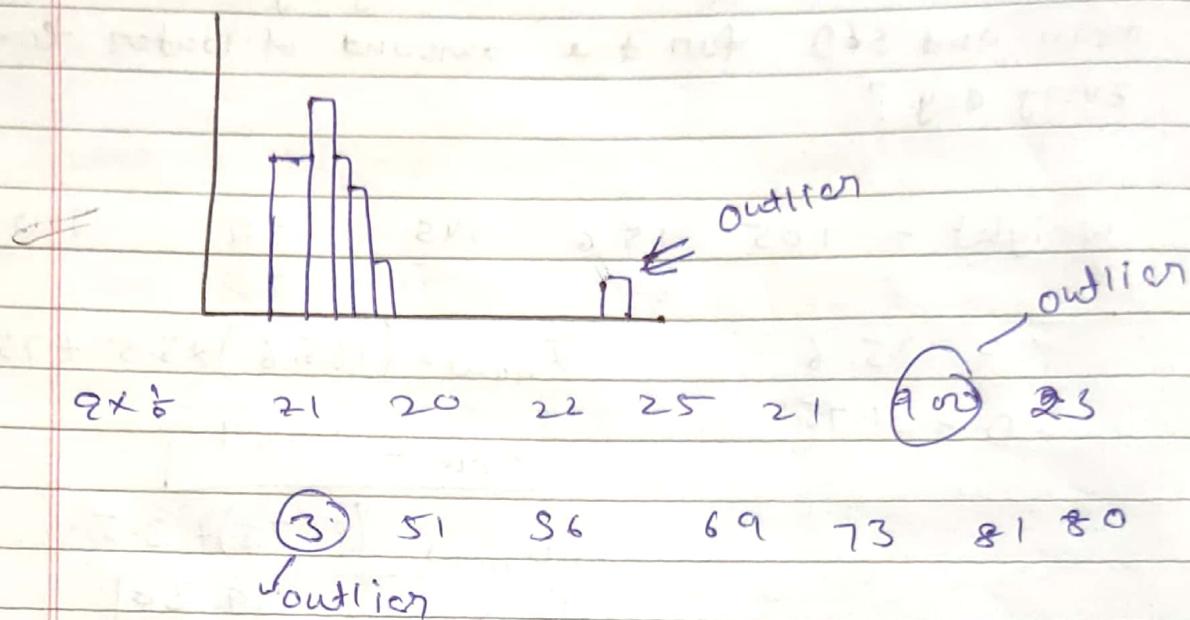
measure of center

$$\boxed{\text{center}_{\text{new}} = (\text{center}_{\text{old}}) * k + b}$$

measure of spread

$$\boxed{\text{spread}_{\text{new}} = (\text{spread}_{\text{old}}) * k}$$

The effect of outlier on spread and center



Note :-

outlier can effect measure of center and spread.

Eg:- temperature of winnipeg on july

Year	Temp	$\bar{x} = -28^{\circ}$	not possible
2015	26°		
2014	15°		
2013	20.5°		
2012	31°	31°	outlier
2011	-350°	-350°	mean is effected by outlier
2010	31°		
2009	30.5°		

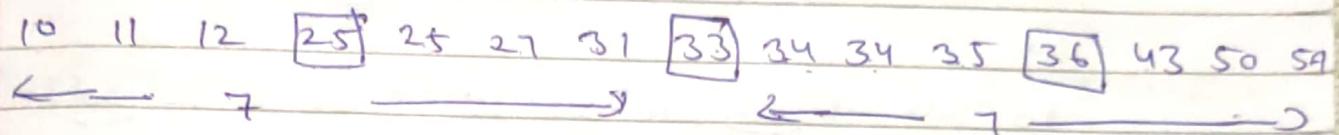
- Range is also effected by outlier
- as we know mean is effected by outlier so as a result S.D also effected by outlier

#

Five Number Summary

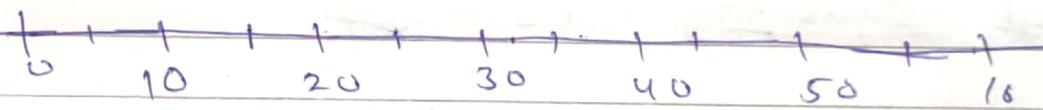
Min	1 st Quartile	Median	3 rd Quartile	Max
10	25	33	36	59

Data Set

⇒ Boxplot :-

Gives us a visual representation of the 5 No. summary

$$\text{IQR} = Q_3 - Q_1 \quad \text{Interquartile Range}$$

⇒ modified Boxplot :-

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

A data value is considered to be an outlier if

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

or

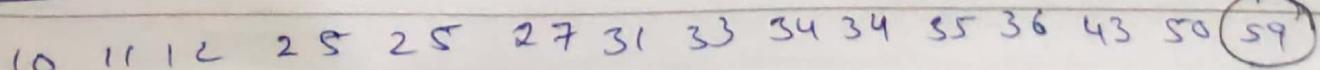
$$\text{Data value} > Q_3 + 1.5(\text{IQR})$$

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

Data value is an outlier if it is

less than : $Q_1 - 1.5(\text{IQR}) = 8.5$

greater than : $Q_3 + 1.5(\text{IQR}) = 52.5$ outlier

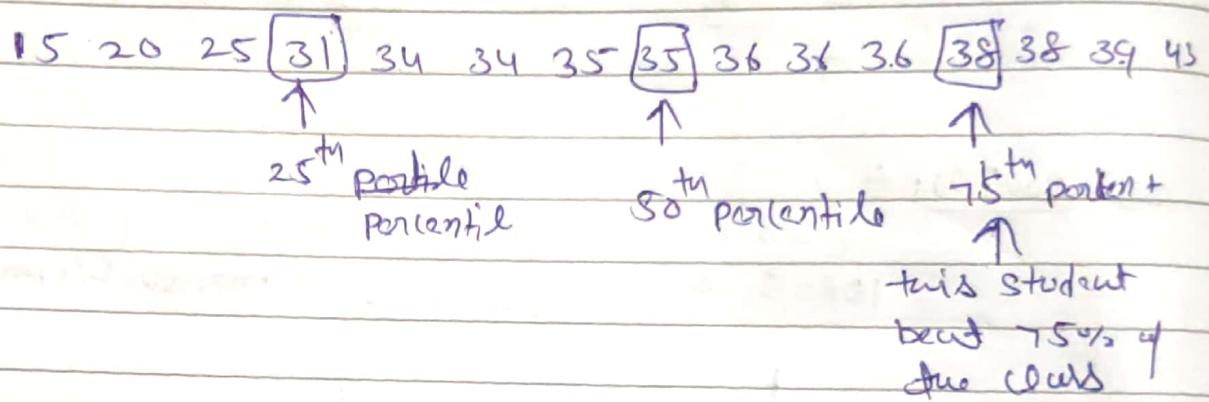
Dataset

#

Percentile :-

Describes the % of data values that fall at or below another data value

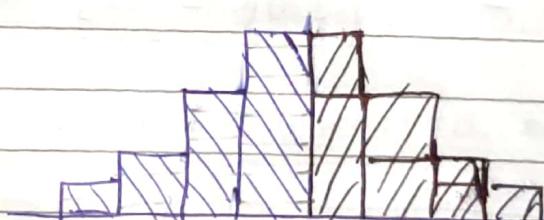
Q Suppose we have 15 student in class and a test is conducted of 50 marks, mark obtain by student is



#

Symmetry and Skewness

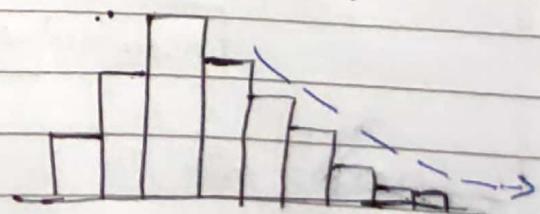
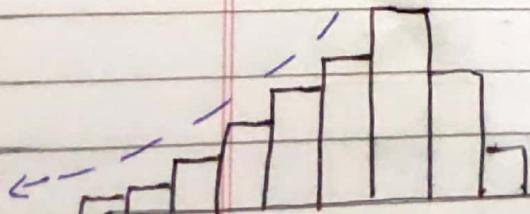
Distribution: Symmetric



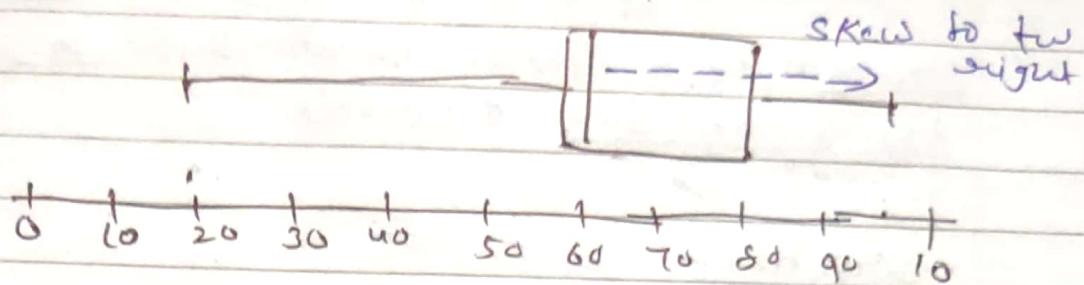
Skew to the left

Skewness Refer to Asymmetry

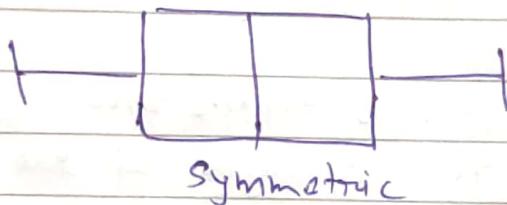
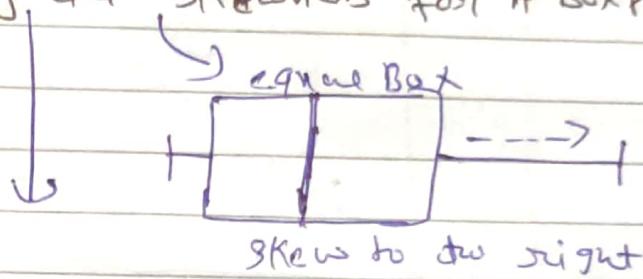
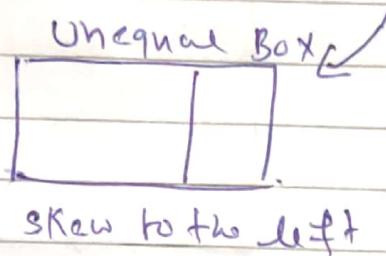
Skew to the right



Box plot = 21, 50, 51, 52, 60, 60, 61, 61, 62, 63, 63, 70, 70
71, 71, 80, 80, 80, 90, 93



Strategies for determining the skewness for a Box plot



Note :-

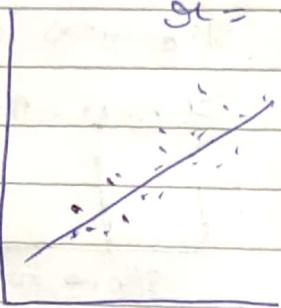
1. if Distribution is Skewed to left then mean is less than median
2. if distribution is skewed to right then median is greater than mean.

correlation :-

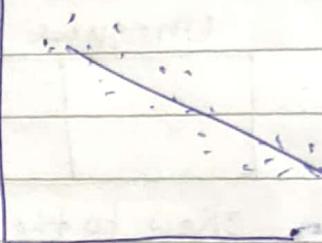
It tells you about the direction and strength of a linear relationship shared b/w two quantitative variables.

Direction

$$\text{or} = +ve$$



$$\text{or} = -ve$$

Strength :-

- or has values b/w 1 and -1
- The strength of the linear relationship increase as or get close to 1 or -1

$$\text{correlation (or)} = \frac{1}{(n-1) S.D_x S.D_y} \sum (x_i - \bar{x})(y_i - \bar{y})$$

Ex:- A teacher wants to determine the correlation b/w. the no. of hours spent studying and test score.

x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$
13	53	0.5	-15	-7.5
15	69	2.5	1	2.5
7	92	-5.5	24	-132
3	10	-9.5	-58	-551
10	85	-2.5	17	-42.5
27	99	14.5	31	449.5

$$\bar{x} = 12.5 \quad \bar{y} = 68$$

$$\underline{821}$$

$$S.D_x = 8.28 \quad S.D_y = 32.91$$

$$g_1 = \frac{1}{(6-1)(8.28)(32.91)} \stackrel{(821)}{=} 0.602$$

$$g_1 = 0.602$$

Regression And R-Square

Regression line:

It predict the change in "y" when "x" increase by one unit

increase/decrease
↓

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

Predicted value at y Intercept slope
Any value of x

$$b_1 = g_1 \frac{S.D_y}{S.D_x}$$

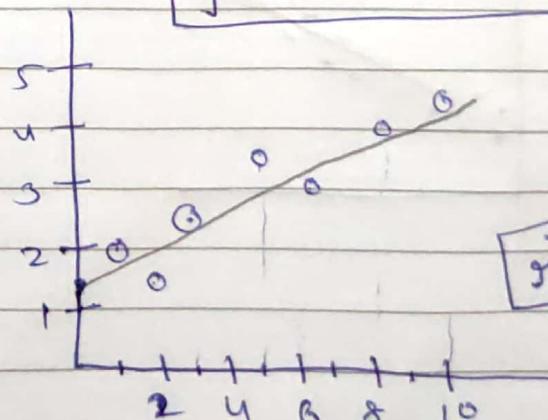
$$b_0 = \bar{y} - b_1 \bar{x}$$

Inq

Q: suppose researcher want to predict a student's GPA from the amount of time they study each week

study time	y_i
1	2.0
2	1.5
3	2.5
5	3.5
6	3.0
8	4.0
10	4.5

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



$$g_1 = 0.88$$

$$\begin{aligned} \bar{x} &= 5 \\ S.D_x &= 3.26 \end{aligned}$$

$$\begin{aligned} \bar{y} &= 3 \\ S.D_y &= 1.08 \end{aligned}$$

$$\begin{aligned} g_1 &= 0.94 \\ b_1 &= 0.311 \end{aligned}$$

$$b_0 = 1.45$$

R-Squared :-

gr

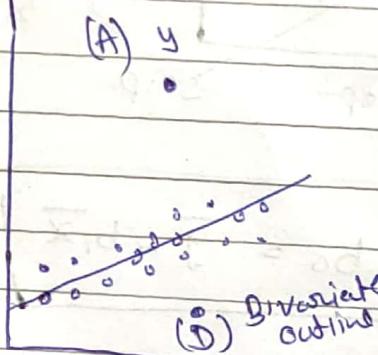
1. Has value $b/w -1$ and 1
2. measures the relationship b/w two quantitative variables with respect to direction and strength

r^2

Has value b/w 0 and 1
is a measure of how close each data point fit to the regression line. it tell us how well the regression line predict the actual values

effect of outlier and extrapolation on Regression

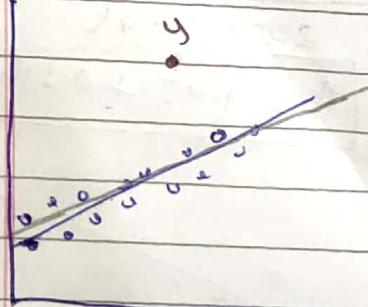
(A) y



(B) (x, y)

without outliers

(C) x



(C) x
influential

bcz it
effect a lot

