

STATISTICS

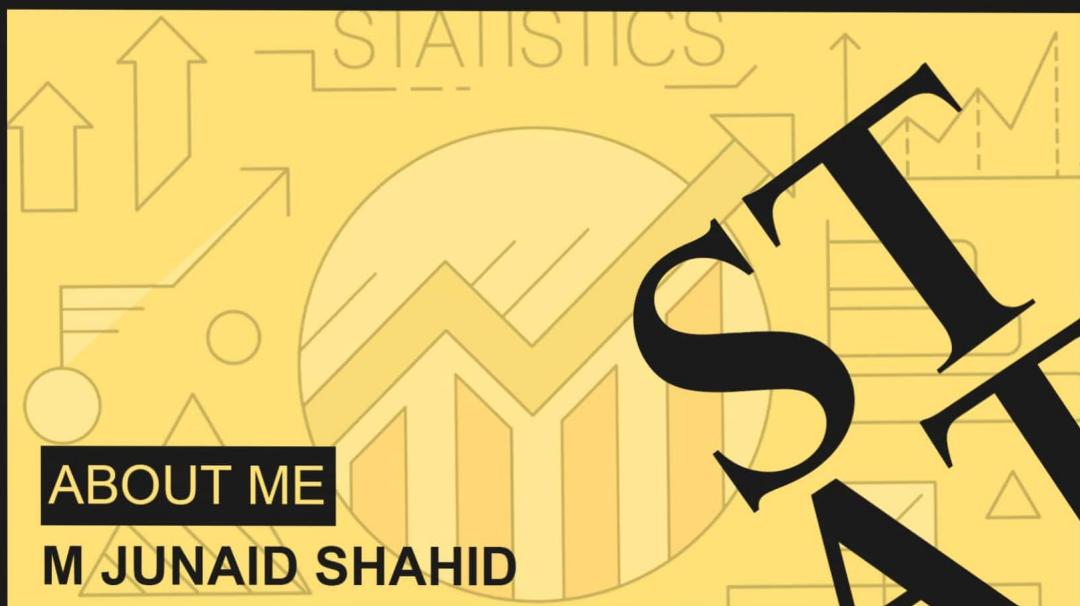
DATA SCIENCE

ST
AT

lecture 1



<https://www.linkedin.com/in/iam-junaid/>



I am a student of MATHEMATICS and studying in the field of data science and data analysis. I am learning online at CAMPUSX DSMP (data science mentorship program). I have done PYTHON with numpy, pandas, matplotlib, seaborn, SQL(structured query language), MATH(linear algebra, vectors, calculus), STATISTICS(Inferential statistics, descriptive statistics, and theorems) and I am still learning to get and polish my new skills to secure my dream DATA SCIENTIST job.

Descriptive Statistics

What is Statistics :-

- Statistics is a branch of mathematics that involves collecting, analysing, interpreting and presenting data. It provides tools and methods to understand and make sense of large amounts of data and to draw conclusions and make decisions based on the data.

In practice, statistics is used in a wide range of fields, such as business, economics, social sciences, medical and engineering. It is used to conduct research studies, analyze market trends, evaluate the effectiveness of treatments and interventions, and make forecasts and predictions.

Example:-

Business (identifying customer behaviour)

Medical (identify efficacy of new medicine (clinical trials))

Government & politics (conducting surveys)

Environmental Science (climate research)

Types of Statistics

↳ Descriptive Statistics

↓

Inferential Statistics

Descriptive Statistics :-

Descriptive

statistics deals with the collecting, organization, analysis interpretation and presentation of data it focus on summarizing & describing the main feature of the set of data. without making inference or prediction about the large population.

Inferential Statistics :-

deals with making conclusion and predictions about a population based on Sample . it involves the use of probability theory to estimate the likelihood of certain event occurring , hypothesis testing to determine if a certain claim about a population is supported by the data and regression analysis to examine the relationship between variables.

Population Vs Sample:-

=>Population:-

population refers to the entire group of individuals or object of that we are interested in Studying it is complete Set of observation that we want to make inference about . For example The population might

be all the students in a particular school or all the cars in a particular city.

Sample :-

A Sample, on the other hand, is a subset of the population. It is a smaller group of individuals or objects that we select from the population to study. Samples are used to estimate characteristics of the population such as the mean or the proportion with a certain attribute. For example, we might randomly select 100 students.

Example :-

- ① All students \rightarrow population
visit school for lectures \rightarrow Sample

Things to be carefull while creating Samples.

- 1 Sample Size
- 2 Random
- 3 Representative.

Parameter Vs Statistics

Parameter:-

A parameter is a characteristics of a population.

Statistics:-

A Statistics is a characteristics of a Sample.

Some of Topic that Come under inferential Statistics

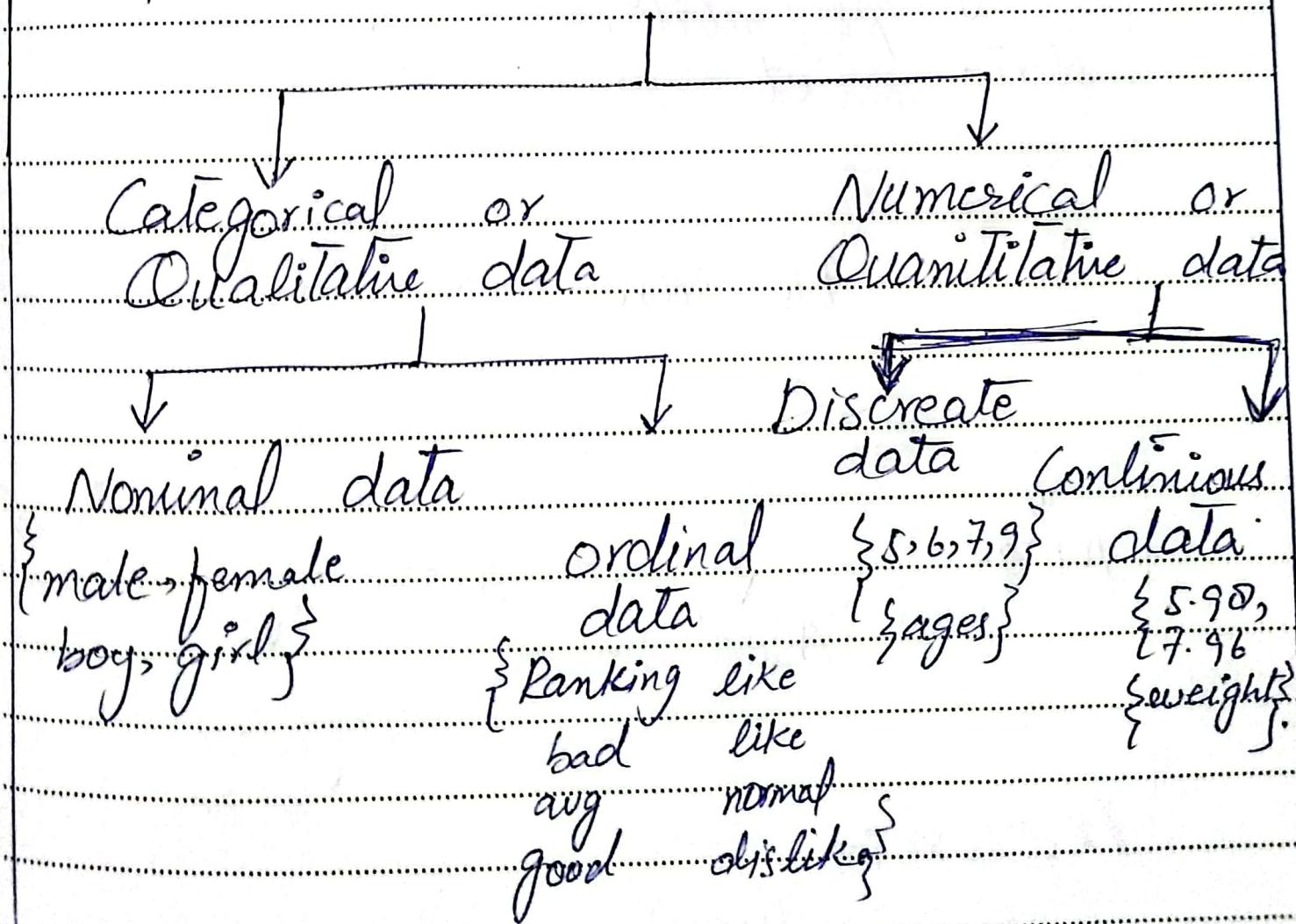
=> Hypothesis Testing:-

=> Confidence Interval:-

- ⇒ Analysis of variance (ANOVA)
- ⇒ Regression Analysis
- ⇒ Chi-Square Test
- ⇒ Sample techniques
- ⇒ Bayesian Statistics

{ Types of Data }

Types of Data



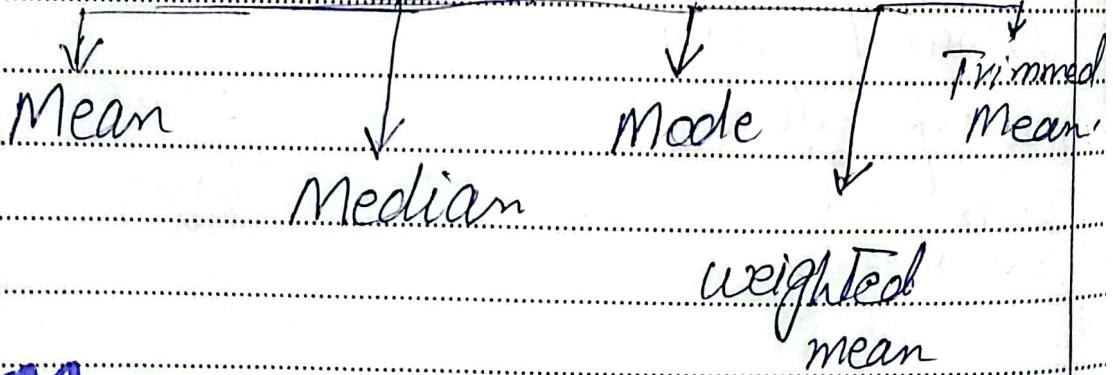
- 3 (a) Mu'awiya was opposed to 'Ali's caliphate. Give reasons for his opposition and write a [1] account of the Battle of Siffin which resulted from this opposition.
- (b) In your opinion what was the most serious consequence of the outcome of this battle? Give [4] reasons for your answer.

Measure of Central Tendency

A measure of Central Tendency is a statical measure that represents a typical Central value of a data set.

It provides a summary of a data by identifying a single value that is most representative of the data set as whole.

Measure of C.T



Mean:-

The mean is a sum of all values in

the dataset divided by the number of values.

Population Mean

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

(N) → no. of items in population

Sample Mean

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

(n) → no. of items in sample.

Marks

5

4

3

2

1

7

$$\bar{x} = \frac{5+4+3+2+1+7}{6}$$

$$\bar{x} = 3.66667$$

Note:-

But if there an outlier in our data set so the mean

will be not accurate

like

Marks with outlier

1

2

3

5

7

9

10

69

$$\bar{x} = \frac{1+3+2+5+7+9+10+69}{8}$$

$$\bar{x} = 13.25$$

$$\bar{x} = 13.25$$

→ outlier

Marks

without outlier.

1

2

3

5

7

9

10

$$\bar{x} = \frac{1+2+3+5+7+9+10}{7}$$

$$\bar{x} = 5.28$$

So that's why we used median when an outlier in our data set.

Median:-

The median is

(a) Islamic teachings revolve around six main Articles of Faith. Write about the following two:

- belief in God, and
- belief in angels.

[10]

(b) Why is the belief in angels important for Muslims?

[4]

is a middle middle value
in our data set when
the data is arranged
in order.

Marks

without Outlier

1

5

7

6

4

3

2

Arrange them

1, 2, 3, 4, 5, 6, 7

↓

middle value

Median = 4

Marks

1

2

3

4

5

6

7

Arrang

1, 2, 3, 4, 5, 6, 7, 8, 10

↓

Median = 5

69

70

This is actually
fine from
Mean &
more accurate

Mode :-

The mode is a value that appears most frequently in the dataset.

Marks

1

5

9

7

5

6

5

4

5

9

7

$$\text{Mode} = 5$$



because 5 contain
more times in
dataset.

Weighted Mean :-

The weighted mean is a sum of the product of each value and its weight, divided by the sum of the weights. It is used to calculate a mean when the value on the

data set have different importance or frequency.

Let we predict on model.

$L \cdot R$	\rightarrow accuracy = 0.2	\rightarrow on this 10 lac
$R \cdot F$	\rightarrow accuracy = 0.3	\rightarrow on this 15 lac
$Xgb00$	\rightarrow accuracy = 0.5	on thi 12

product weight value

$$W_m = \frac{\text{Sum of product of value with its weight}}{\text{Sum of weights}}$$

$$= \frac{0.2 \times 10 + 0.3 \times 15 + 0.5 \times 12}{0.2 + 0.3 + 0.5}$$

$$\boxed{W_m = \text{Ans}}$$

Trimmed Mean:-

A trim mean is calculated by removing a certain percentage of the smallest and largest value from the dataset and then taking the mean.

- 5 (a) 'Prophets played a central part in conveying God's message to humanity.' Write an account of Muslim belief in prophets. [1]

- (b) Why do you think God gave miracles to his chosen prophets?

of the remaining values.
the percentage of the
value removed is called
trimmed percentage.

→ without outliers & not trimmed
 $20k, 22k, 23k, 25k, 28k, 30k, 32k, 35k,$
 $50k, 80k$

$\boxed{\bar{X} = 36k}$
→ By trimmed 20% from largest
& smallest side

$25k, 28k, 30k, 32k, 35k$

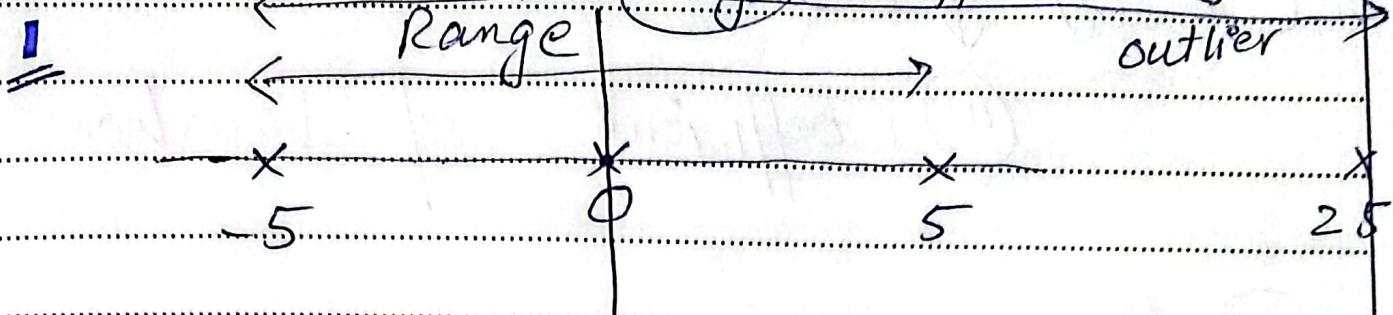
$$\boxed{\bar{X} = 30k}$$

Measure of Dispersion:-

A measure of dispersion is a statistical measure that describes the spread or variability of dataset. It provides information about how much data set is distributed around the

Central tendency (mean, median, mode) of data set

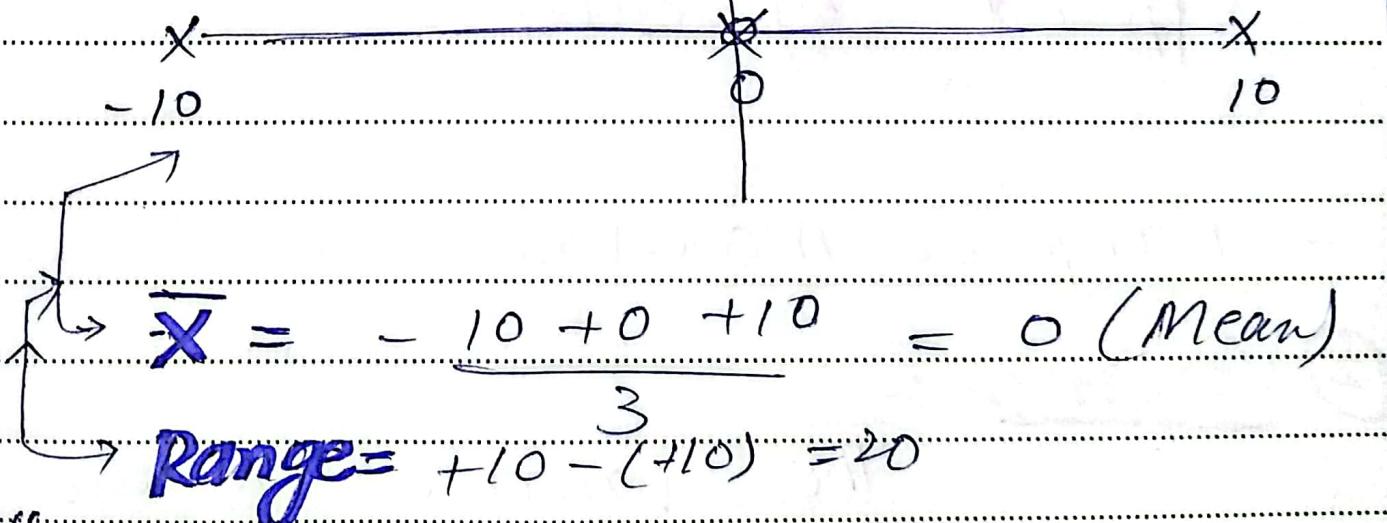
(Range) → effected by



$$\bar{x} = \frac{-5 + 0 + 5}{3} = \frac{0}{3} \Rightarrow 0 \text{ (Mean)}$$

$$\text{Range} = 5 - (-5) = 10$$

2



Both same in two example

So how we identify
identify it by how the
data is spread by measure
of dispersion:

Measure of dispersion

① Range ② Variance ③ S.D

④ Coefficient of Variation

① Range :-

The Range is the difference between the maximum and minimum values in the data set. It is a simple measure of dispersion that is easy to calculate but can be affected by outliers.

$$\Rightarrow \text{Range} = \text{Max-Value} - \text{Min-Value}$$

② Variance :-

The variance is the average of the squared difference between each data

Permission to reproduce items where third-party owned material protected by copyright is included has been sought and cleared where possible. Reasonable effort has been made by the publisher (UCLES) to trace copyright holders, but if any items requiring clearance have unwittingly been included, publisher will be pleased to make amends at the earliest possible opportunity.

To avoid the issue of disclosure of answer-related information to candidates, all copyright acknowledgements are reproduced online in the Assessment International Education Copyright Acknowledgements Booklet. This is produced for each series of examinations and is freely available online at www.cambridgeinternational.org after the live examination series.

point and the mean. it measures the average distance of each data point from the mean and is useful in comparing the dispersion of dataset with different means.

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N} \text{ population}$$

$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \text{ Sample}$$

Data $x_i - \bar{x}$ $(x_i - \bar{x})^2$

$$3 \quad 3-3 \quad 0 \quad X = \underline{3+2+1+5+4}$$

$$2 \quad 2-3 \quad 1 \quad 5$$

$$1 \quad 1-3 \quad 4 \quad = 3$$

$$5 \quad 5-3 \quad 4$$

$$4 \quad 4-3 \quad 1$$

$$S^2 = \frac{0+1+4+4+1}{5}$$

$$S^2 / 6^2 = \boxed{\text{Ans}}$$

11

Note

So we used Square on each

Value for if we get the answer in negative after to sort out we used square one person said we can used absolute in / or other hand \Rightarrow Yes, we can but one markdown is the absolute formula not provide inference about data we will use the name of Question is

\Rightarrow Mean absolute deviation :-

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

STANDARD DEVIATION:

is a square root of variance.
it is widely used measure
of dispersion that is useful
in describing the shape
of distribution.

its also used to change unit of data same as data \Rightarrow mean if I have

writing?

If i have a dataset have unit LPS then by using variance the unit is $(\text{LPS})^2$ & again taking S.D it convert the $(\text{LPS})^2$ it in same unit off data.

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{N}} \text{ Population}$$

$$S = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}} \text{ Sample}$$

→ Coefficient of Variation:-

Coefficient of variation (CV) The ratio of the standard deviation to the mean expressed as a percentage. It is used to compare the variability of datasets.

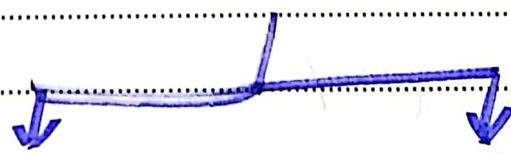
with different mean and
is commonly used in
field such as biology,
chemistry & engineering.

Formula

$$CV = \frac{\text{Standard deviation}}{\text{Mean}} \times 100$$

it tells us also the
distribution between two
Comparing Quantities how
much data is spread.

☰ Graphs For univariate Analysis



Categorical Numerical

Categorical :- frequency distribution

A frequency distribution
table is a table that
Summarizes the numbers
of times (for frequency)

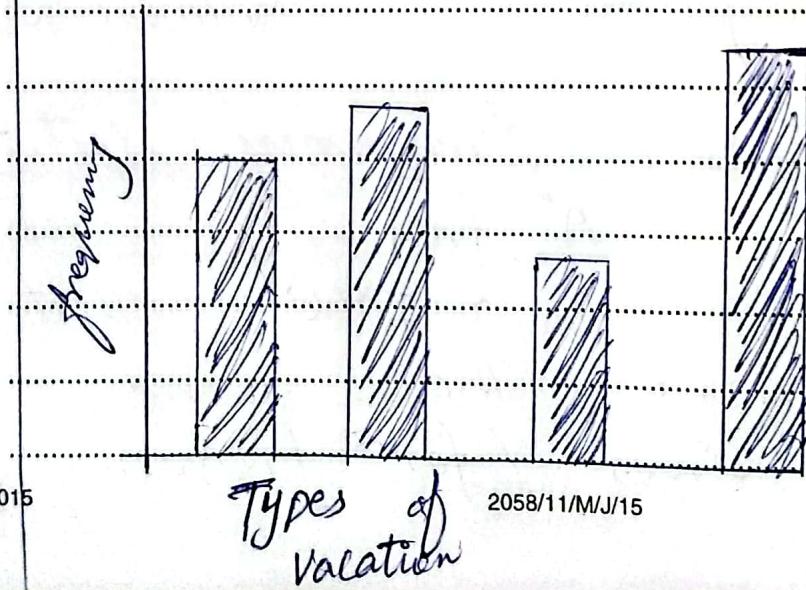
that each value occur
in data set

E.G

let Survey of 200 peoples
and ask about
favourite type of vacation

Type of Vacation	frequency
Beach	60
City	40
Adventure	30
Nature	20
Cruise	15
Other	35

In this we created
bar graph



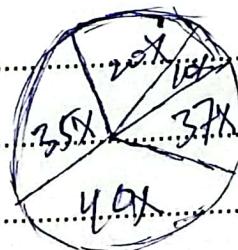
- (a) Write about the changes in the relationship between the Prophet and the Quraysh in the years between his marriage to Khadija and the death of Abu Talib. [10]
- (b) Why is it significant that the Quraysh were still willing to keep their belongings with the Prophet after he started to preach Islam? [4]

Relative frequency: is the proportion or percentage of a category in a data set or sample. It is calculated by dividing the frequency of a category by the total number of observation in the dataset.

Types of vacation frequency $R.F = \frac{freq}{200} \times 100$

	60	=	
Beach	60		0.3
City	40		0.2
Adventure	30		0.15
Nature	35		0.175
Cruise	20		0.1
Other	15		0.075

Pie chart:



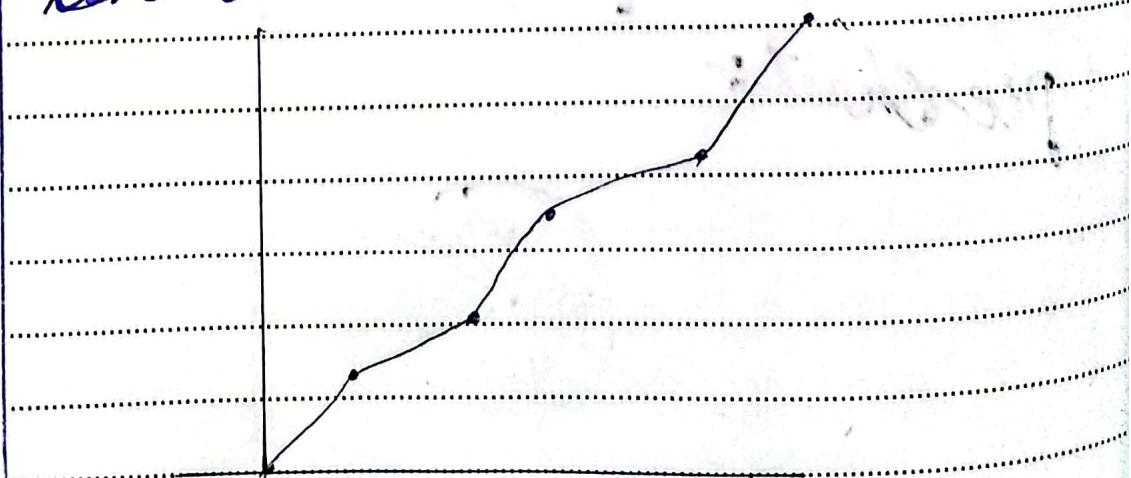
Cumulative frequency:-

is the running total of frequencies of a variable or category in a dataset or sample. it is calculated by adding up the frequencies of the current category and all previous categories in the dataset or sample.

Types of Vacation frequency Cumulative frequency

Beach	60	60
City	40	100
Adventure	30	130
Nature	35	165
Cruise	20	185
other	15	200

line chart:-



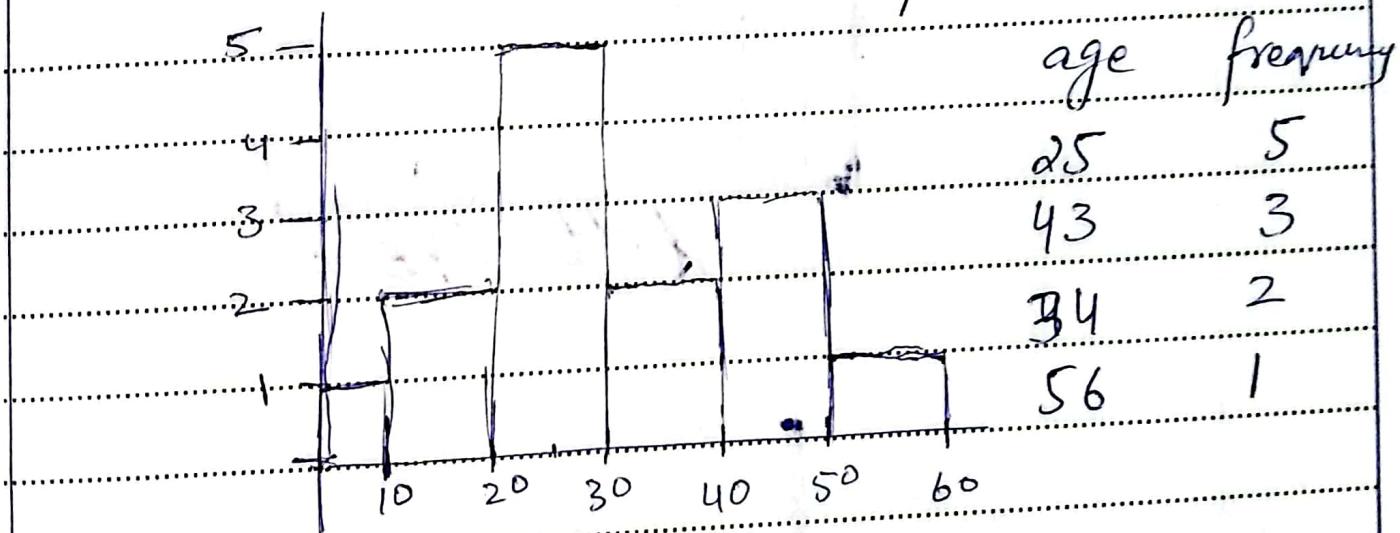
② Numerical:-

In Univariate

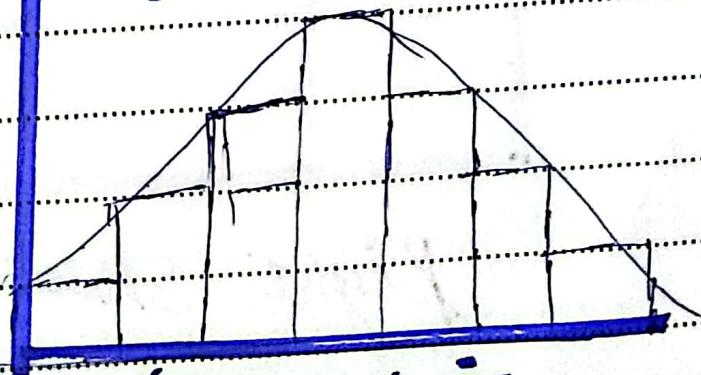
numerical column we can draw frequency distribution graph as well as Histogram.

Histogram:-

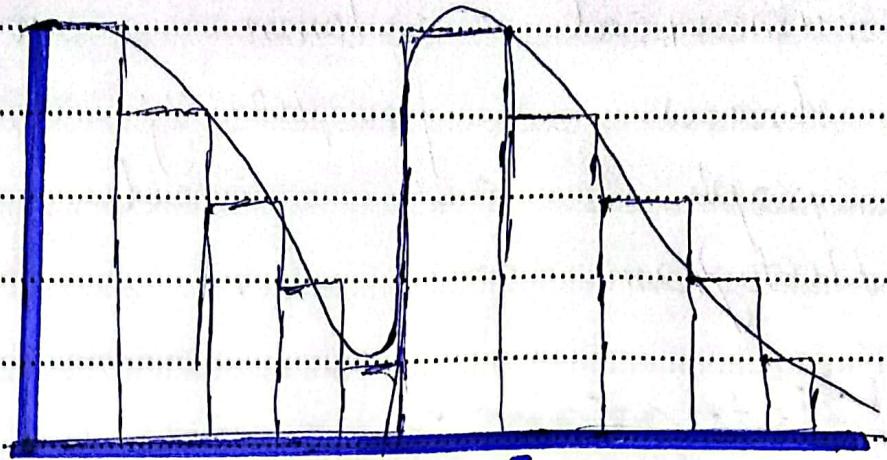
it is a representation of data in a hist plot
In this graph we distribute our data in bins/buckets.



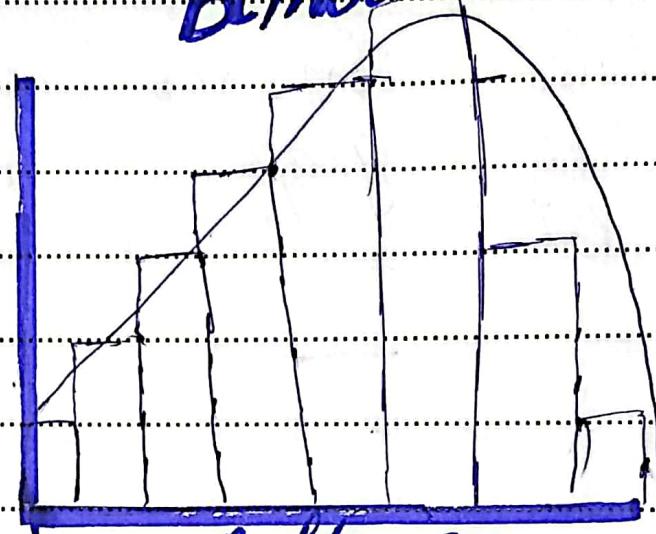
Types of histogram:-



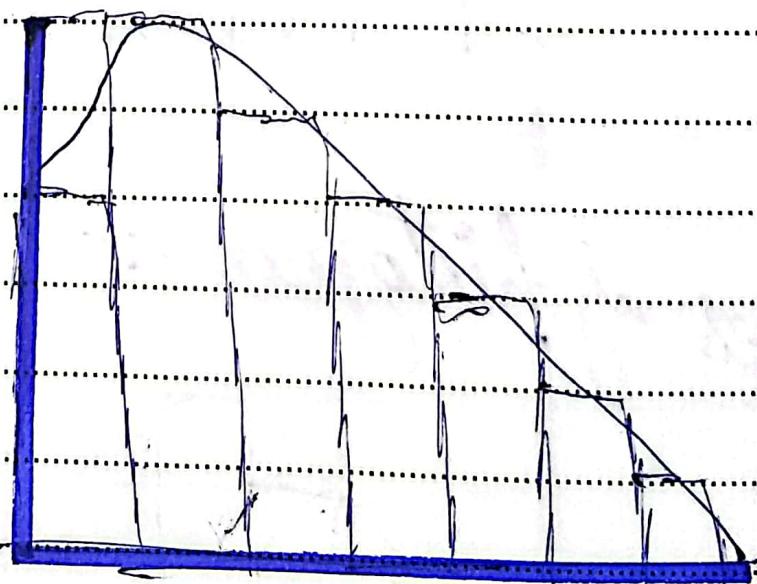
Symmetric



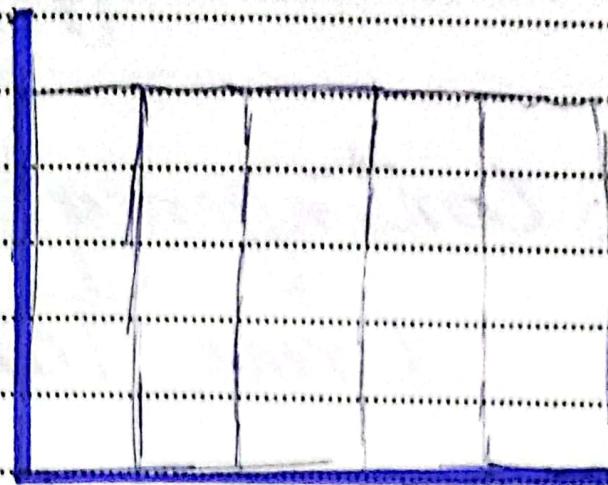
Bimodel



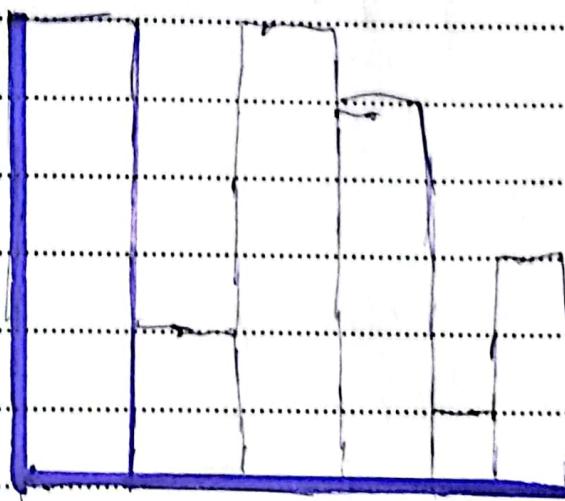
left skew



right skew



Uniform



No-pattern

≡ Graph of Bivariant Analysis

Categorical
Categorical
Categorical

Categorical
Numerical

Numerical
Numerical
Numerical

= Categorical - Categorical

= Consists Contingency Table /
Cross Table

A Contingency Table also known
as cross-Tabulation is
a type of table used in
statistics to summarize the
relationship between two
categorical variables.

A Contingency Table displays
the frequency or relative
frequency of the observed
values of the two
variables organized into
rows and columns.

Survived		Pclass
0	1	1
		2
1	3	3

(b) What lessons can Muslims learn?

PClass

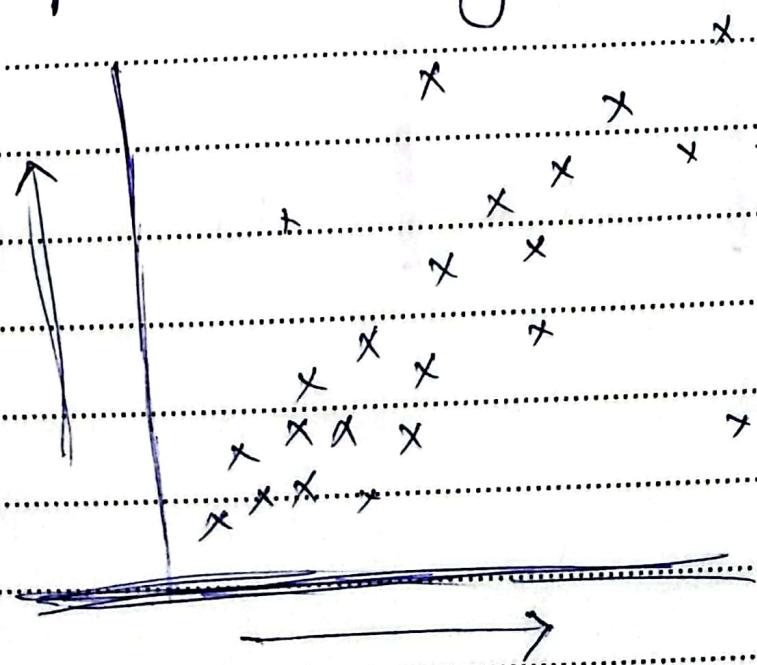
Survived	1	2	3
1	71	31	63
0	42	118	13

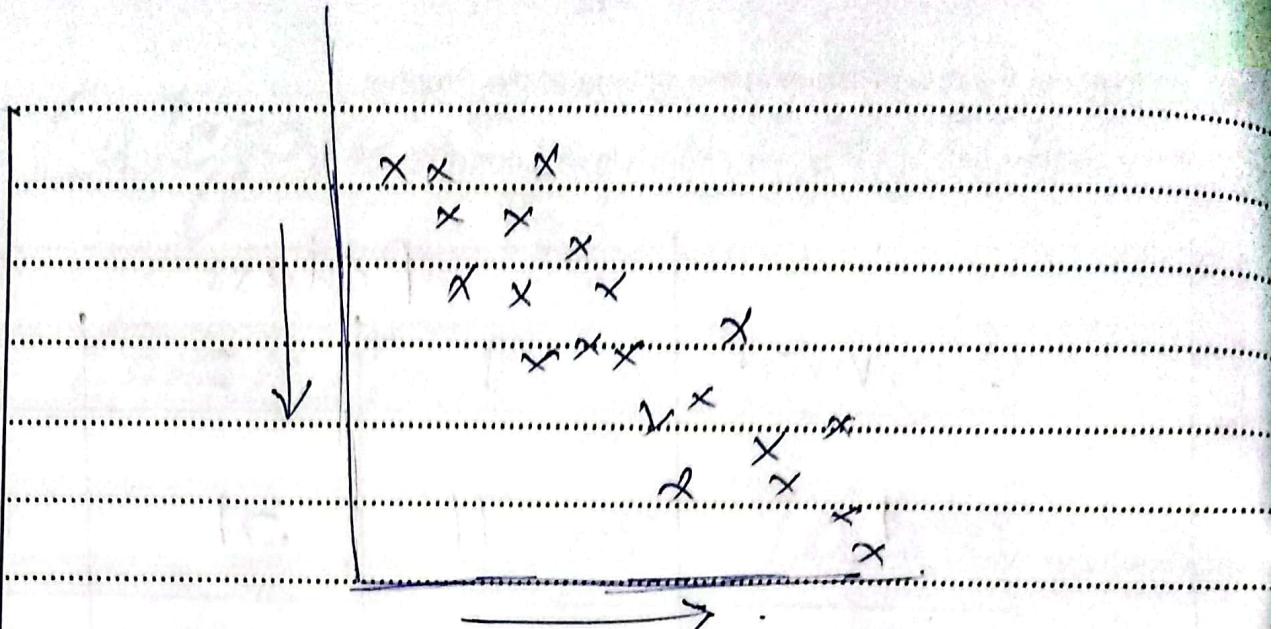
Numerical - Numerical:-

In Numerical - Numerical bivariate analysis we can draw a scatter plot.

point where distribution of data spread using points

Scatter plot:-





Categorical - Numerical

\therefore In Categorical - Numerical
Bivariate analysis
we can make histogram
Contingency table.

	0 - 10	11 - 20	21 - 30
Male	32	41	110
Female	15	18	120

THANKS