

Statistics.

M	T	W	T	F	S
Page No.:	YOUVA				
Date:	YOUVA				

- * * What is statistics and its types :-
- * Definition :- Statistics is the science of collecting organizing and analyzing data and interpreting data to support decision making and understanding.
- * Data :- Facts or pieces of information.
- * Types of statistics :-

Descriptive statistics

(understanding data)

Defn:- It consist of organizing and summarizing of data.

a) Measure of central tendency

(mean, mode, median)

b) Measure of dispersion

(variance & standard deviation)

c) Shape

Different type of distribution of data e.g. skewness, kurtosis.

d) Histogram, probability density function & frequently probability mass function. (Visually)

of the element.

Inferential statistics

(Drawing conclusions of data)

It consist of using data; you have measured to form conclusion.

a) Z-test

b) T-Test

c) Anova test

d) Chi-square test

} hypothesis testing

Sample data

Population

for descriptive statistics.

What → summarization

Why → find pattern &ings

When → before decision

making

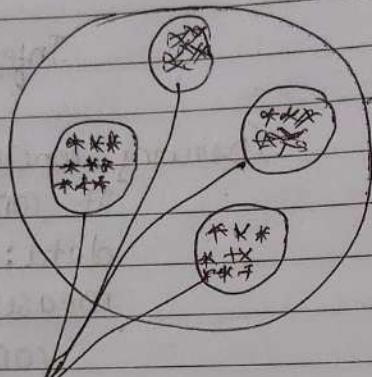
Where → everywhere.

* Types of Data :-

* Sampling Technique

① population - The group that we are interested in studying.

② Sample - A subset of the population.



Sample.

The goal of sampling is to create a sample that is representing of the entire population.

$$\text{Population} = (N)$$

$$\text{Sample} = (n)$$

③ Types of sampling :-

1) Simple random sampling :- When performing simple random sampling every member of the population (N) has an equal chance of being selected. for your sample (n)

2) stratified sampling

Layering (Non overlapping groups)

population
 male female
 They do not overlap.

2) Systematic sampling :- mall \rightarrow outside

n^{th} person 4^{th} person
 will choose for ↓
 survey. survey.

4) Convenience sampling :-

Voluntary Response

Data Science
 Survey

sampling

Data science knowledge.

* What are variables & it's types :-

Defn - Variable is a property that can take on many values.

e.g. $\text{age} = 12$

↓
 variable

13

40

100

$\text{height} = 172 \text{ cm}$

↓
 172.5 cm

variable

$\text{weight} = 72 \text{ kg}$

↓
 72.5 kg

↓
 73 kg

variable 100 kg.

Variable is only in singular mode

Ages = { 15, 20, 25, 30 }

plural = that's why it's not variable.

* Types of Variables :-

Quantitative variable

- Discrete variable (continuous variable.)

- eg. whole number

& no decimal value

especially in output.

eg 1) No. of bank account

2) No. of children.

eg. 1) height =

175.50 cm

2) weight =

72.3 kg

Qualitative or
(categorical) variable

properties

Classification

eg. 1) Gender male
female

2) Type of flower

Rose lotus lily

Nominal Ordinal

(No order matter) (order matters)

* Measure of Central Tendency :-

- ① Mean
- ② Median
- ③ Mode

Defn

Central tendency refers to the measure used to determine the "center" of the distribution of data.

$$\{1, 2, 2, 3, 4, 5\}$$

① Mean (Average)

Population Data (N)

$\mu = \text{population mean}$

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

Sample data (n)

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

Σ = summation

x_i = particular data point

i = row number.

$$1+2+2+3+4+5$$

$$\mu = \frac{1+2+2+3+4+5}{6}$$

$$\mu = \frac{17}{6}$$

$$\mu = 2.83$$

If we will focus on first mean i.e. 2.83 & after adding one outlier it turned to 16.71. But when we worked on median even though we added two outliers it's 3.5. This is how median deal with outliers.

② Median $\circ - x = \{1, 2, 2, 3, 4, 5, 100\}$

help to deal with outliers

mean (avg) = 16.71

$$\mu = \frac{1+2+2+3+4+5+100}{7}$$

A datapoint who doesn't follow the pattern of data or are extreme points.

$$\mu = \frac{117}{7}$$

$$\text{mean} = \mu = 16.71$$

* Sort all the numbers

* find the central element. $\begin{cases} \text{odd length} \\ \text{even length} \end{cases}$

$$\{1, 2, 2, 3, 4, 5, 100\} \rightarrow \text{odd length } \left(\frac{n}{2} \right)^{\text{th}} + \left(\left(\frac{n}{2} \right) + 1 \right)^{\text{th}}$$

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}}$$

$$\{1, 2, 2, 3, 4, 5, 100, 101\} \rightarrow \text{even length. } \left(\frac{n}{2} \right)^{\text{th}}$$

$$\frac{3+4}{2} = 3.5 \quad \text{Median}$$

③ Mode :- Most frequent element. We mostly use of categorical feature.
 $\{1, 2, 2, 3, 3, 3, 4, 4, 4, 5, 5, 5\}$

$$\text{Mode} = l + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h$$

Age question :- find mode of following data:-

Age (in years)	No. of patients (f)
5 - 15	6
15 - 25	11
25 - 35	21
35 - 45	(23) 21
45 - 55	14
55 - 65	5

difference b/w $h = 15 - 5$
 class interval $= 10$
 (or range)

$$\text{Mode} = l + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h$$

$$l = 35$$

$$f_0 = 21$$

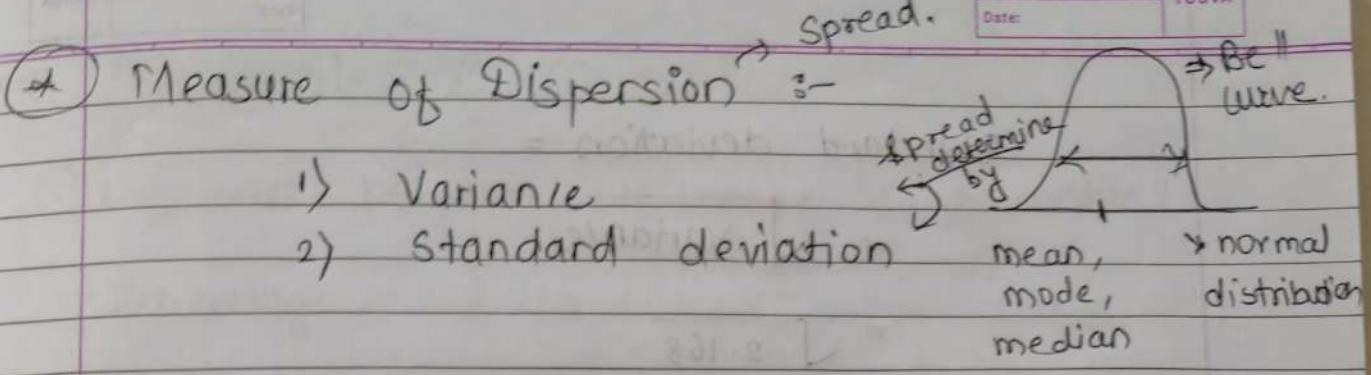
$$f_2 = 14$$

$$h = 10$$

$$f_1 = 23$$

$$= 35 + \left(\frac{23 - 21}{2(23) - 21 - 14} \right) \times 10$$

$$= \underline{\underline{36.8}}$$



* Mean :-

sample mean
(n)

population mean
(N)

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

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

* Variance :-

Sample variance

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

Bessel correction

or also known as
degree of freedom

population variance

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

(population variance)

e.g. $x = \{1, 2, 3, 4, 5\}$

x	\bar{x}	$(x - \bar{x})$	$(x - \bar{x})^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
3	2.83	-0.83	0.6889
4	2.83	0.17	0.03
5	2.83	1.17	0.137
			4.71
$\bar{x} = 2.83$			10.84

$$S^2 = \frac{10.84}{5} =$$

sample variance \Rightarrow {spread of the data}

Standard deviation =

$\sqrt{\text{variance}}$

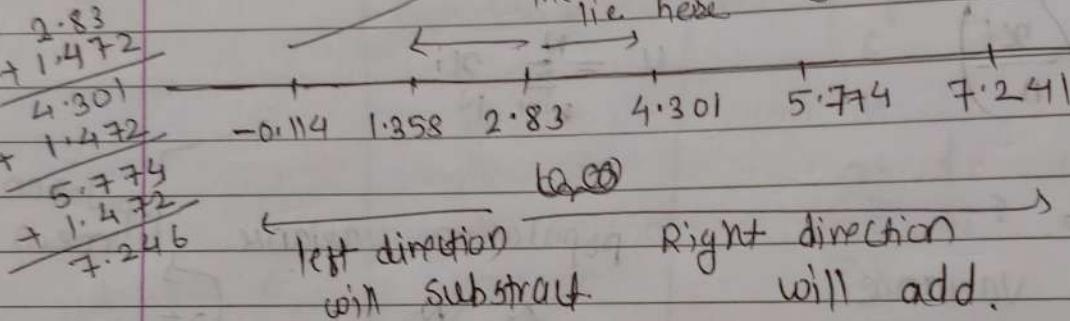
$$= \sqrt{2.168}$$

Sample

standard deviation = 1.472

$$\bar{x} = 2.83$$

$$S = 1.472$$

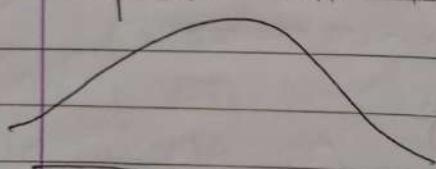


Variance decide the spread.

Variance \rightarrow Big number

Sid. \rightarrow Big number

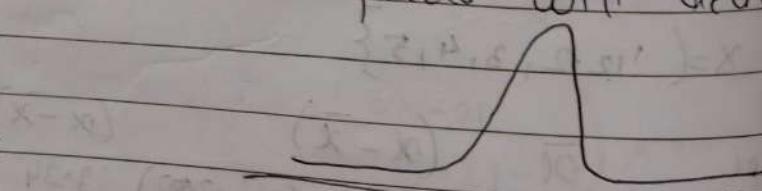
Spread will increase



Variance \rightarrow small number

S.d. \rightarrow small number

spread will decrease



** Percentile & Quartiles :-

(*) Difference between percentage & percentile :-

Percentage :- {1, 2, 3, 4, 5}

% of the numbers that are even ?

$$\% \text{ of even} = \frac{2}{5} = 0.4 = 40\% \text{ of the numbers are even.}$$

Percentiles :- {CAT, GATE, SAT}

Defn :- A percentile is a value below which a certain percentage of observations lie.

95 percentile means that the person has got better marks than 95% of the entire student.

e.g.

(1) Dataset :- {2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12}

$$n = 20$$

① What is the percentile ranking of 10?

$$\text{percentile rank of } 10 = \frac{\# \text{ of values below } 10}{n} \times 100$$

$$= \frac{16}{20} \times 100$$

= 80% (percentile)

② What is the percentile ranking of 11?

③ What value exist at percentile rank of 25?

$$\text{value} = \frac{\text{Percentile}}{100} \times (n+1)$$

$$= \frac{25}{100} \times 21 = 5.25 \quad (\text{This is the Index value})$$

now at 5 position there is number 5 and at here its 5.25 so, will consider next index as well so, it will be 5 again will do the average of both as $\frac{5+5}{2} = \underline{\underline{5}}$
so, the value will be 5



Quartiles :- That means $1/4^{\text{th}}$

~~maximum~~ minimum

Also known as 5-number summary
25% \rightarrow 1st quartile (Q_1)

median

75% \rightarrow 3rd quartile (Q_3)

maximum

Interquartile range (IQR) = $Q_3 - Q_1$
This help to find outliers.

* Scales of Measurement :-

It defines how the data is classified, compared & used mathematically.

* Properties

measurement level

qualitative

Quantitative

- (a) Nominal
- (b) Ordinal

(a) Discrete

(b) Continuous.

* Def'n of Measurement :- Assignment of numbers to characteristics of objects or events according to rules.

* Properties of scales of measurement :-

- ① Identity :- Each value has a unique meaning or ^
- ② Order :- Values can be ranked or ordered logically
- ③ Equal Intervals :- The difference between values is meaningful and equal.
- ④ True zero :- A value of zero means the absence of the quality. (e.g. zero weight = no weight.)

~~Types of
Data from
pdf~~

Nominal :- categorizes data without a specific order. e.g. Gender, blood type

ordinal :- categorizes data with a meaningful order. e.g. education level, ranking,

Interval :- Numeric data with ordered categories and equal intervals but no true zero. e.g. Temp., IQ score

Ratio :- Like interval but with a true zero point. e.g. Height, weight, age, income

* Interquartile Range (IQR)

$$IQR = Q_3 - Q_1$$

$$Q_3 = 9 \quad Q_1 = 5$$

$$IQR = 4$$

~~rough example~~ min \rightarrow lower fence = $Q_1 + 1.5(IQR)$
 $= 5 + 1.5(4)$
 $= 5 + 6$
 $= -1$

max \rightarrow higher fence = $Q_3 + 1.5(IQR)$
 $= 9 + 1.5(4)$
 $= 9 + 6$
 $= 15$

$x < -1$ = outlier

$x > 15$ = outlier.

Box plot min = -1

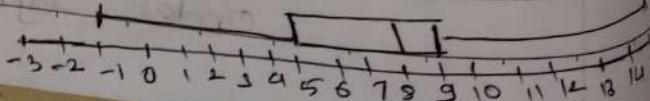
25th = 5

median = 8

75th = 9

Max = 15

$Q_3 - Q_1$



2nd example :-

$$x = \{ 1, 2, 2, 2, 3, \underline{3}, 4, 5, 5, 5, 6, 6, 6, 7, 8, \underline{8}, 9, 15, 12, 7 \}$$

$$Q_1 = \left(\frac{25}{100} \times 20 \right) + 1$$

= 6th Index

= On 6th Index = 3

$$Q_3 = \left(\frac{75}{100} \times 20 \right) + 1 = 16$$

= 16th Index

= 8

$$IQR = Q_3 - Q_1$$

$$= 8 - 3$$

$$= 5$$

$$\begin{aligned} \text{Min-Lower fence} &= Q_1 - 1.5(IQR) \\ &= 3 - 1.5(5) \\ &= -4.5 \end{aligned}$$

$$\begin{aligned} \text{Max-Higher fence} &= Q_3 + 1.5(IQR) \\ &= 8 + 7.5 \\ &= 15.5 \end{aligned}$$

$$\text{min} = -4.5$$

$$Q_1 = 3$$

$$\text{median} = 5.5$$

$$Q_3 = 8$$

$$\text{max} = 15.5$$

* Descriptive statistics. *

Defn

D.S. are the methods descriptive statistics is a branch of statistics that summarizes and describes the main feature of a dataset. It helps us understand the basic characteristics of data through numerical measures and graphical representations.

(A)

Types of D.S. :-

1) Measure of central tendency :-

This refers to the measure used to determine the "center" of distribution of data.

(a) Mean :- The average of the data.

(b) Median :- The middle value when data is ordered.

(c) Mode :- The most frequent value.

2) Measures of Dispersion (spread) :-

These show how spread out the data.

(a) Range :- Difference between the maximum & minimum values.

(b) Variance :- Average of the squared differences from the mean.

(c) Standard deviation :- Square root of the variance shows average distance from the mean.

- ④ Interquartile Range (IQR) :-
Difference between the 75th percentile (Q_3) and 25th percentile (Q_1)
- 3) shape of the Distribution
- a) skewness :- Measures the asymmetry of the distribution.
- b) kurtosis :- Measures the "tailedness" of the distribution.

(B) Data Visualization Tools :-

- 1) Histogram :- shows the distribution of data.
- 2) Bar chart :- Useful for categorical data.
- 3) Pie chart :- Shows proportions of categories

Descriptive Statistics

Measure of
Central Tendency

- 1) Mean
- 2) Median
- 3) Mode

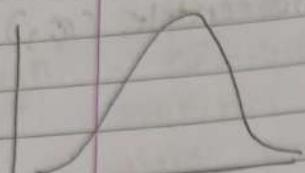
Measure of
Variability / Dispersion

- 1) Range
- 2) IQR
- 3) Variance
- 4) std. deviation

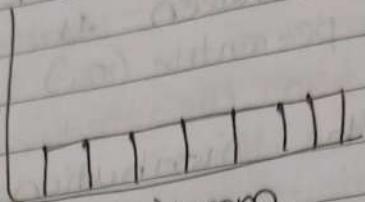
Shape /
freq. distribution

- 1) Skewness
- 2) Kurtosis.

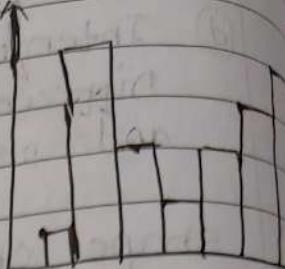
* Different shapes of Distribution :-



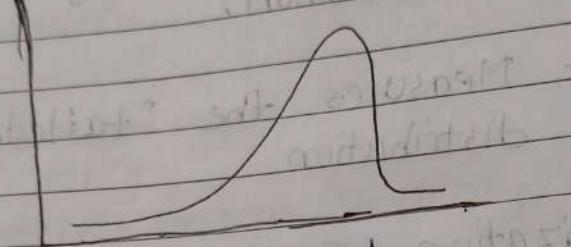
symmetric



Uniform

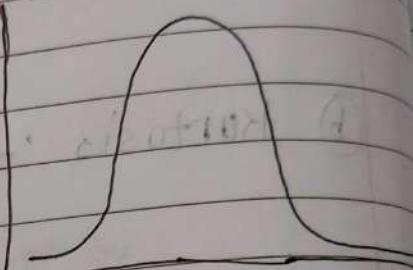


Bimodal



skewed left.

Mean > Median
(+vely skewed)



skewed right
Mean < Median
(-vely skewed)

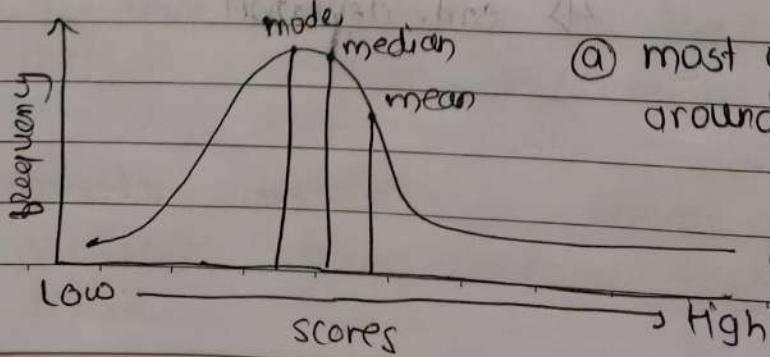
(A) Skewness :- Measures lack of symmetry in data distribution.

(less symmetric data, more the skewness)

* for normal distribution skewness is zero.

* Any kind of symmetrical data shows skewness close to zero.

(a) Right-skewed distribution (+ve skewness)



(a) most of the values are around left of the distribution.

(b) Mean, median greater than the mode.

- (b) Left-skewed distribution
(Negative skewness)

① most of the values are clustered around eight.

② Mean, median less than the mode.

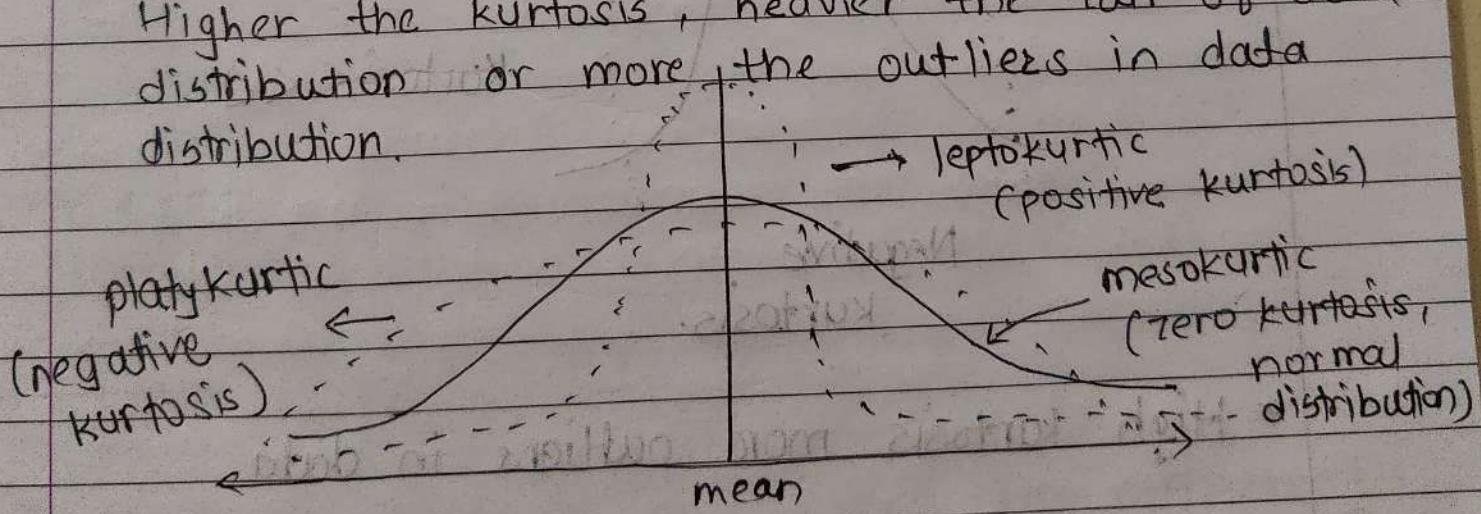
* If skewness is less than -1 or greater than 1 :- Highly skewed.

* If skewness is between -1 and -0.5 or between 0.5 and 1 :- Moderately skewed

* If skewness is in between -0.5 and 0.5 :- Approximately symmetric.

- (C) Kurtosis :- measures whether data is heavy tailed or light-tailed.

Higher the kurtosis, heavier the tail of data distribution or more the outliers in data distribution.



kurtosis distribution

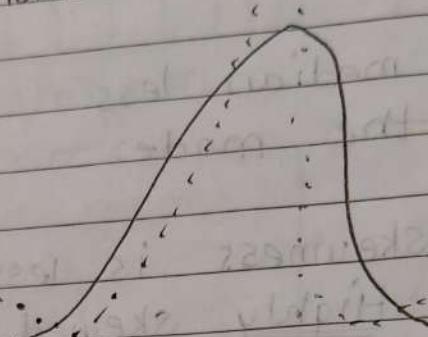
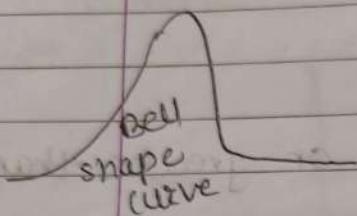
- * negative (-ve)
- * Zero
- * Positive (+ve)

Date: _____
Leptokurtic
(+ve kurtosis)

Mesokurtic
(zero kurtosis)

* looks similar to
the normal distribution
curve.

- * longer distribution
- * sharper peak & heavier tails
- * More outliers.

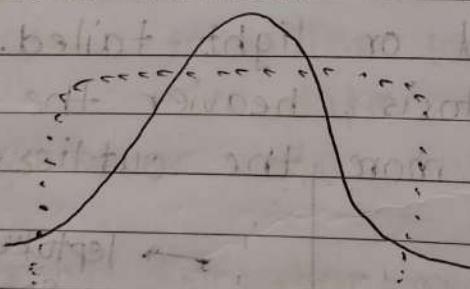


Probability of having
outliers or extreme
values is zero or close
to zero.
Kurtosis = 0

Positive kurtosis

Kurtosis > 3

- * shorter distribution
- * flatter peak & thinner tails
- * Less outliers.



Kurtosis < 3

Negative
kurtosis.

- * High kurtosis more outliers in data.

④ When is kurtosis useful :-

- Only when it's used in conjunction with standard deviation.

- It is possible that a random variable might have a high kurtosis (bad) but the overall standard deviation is low (good).
- Conversely, one might see another random variable with a low kurtosis (good), but the overall standard deviation is high (bad).
- * Does sample size have an impact on skewness & kurtosis?
They are always increased with increasing sample size.

(A) Skewness :-

$$\frac{1}{(n-1)(n-2)} \leq \left[\frac{(x_i - \bar{x})^3}{S.D.} \right]$$

x_i = data point
 \bar{x} = mean
 $S.D.$ = standard deviation
 n = no. of observation

(B) Kurtosis :-

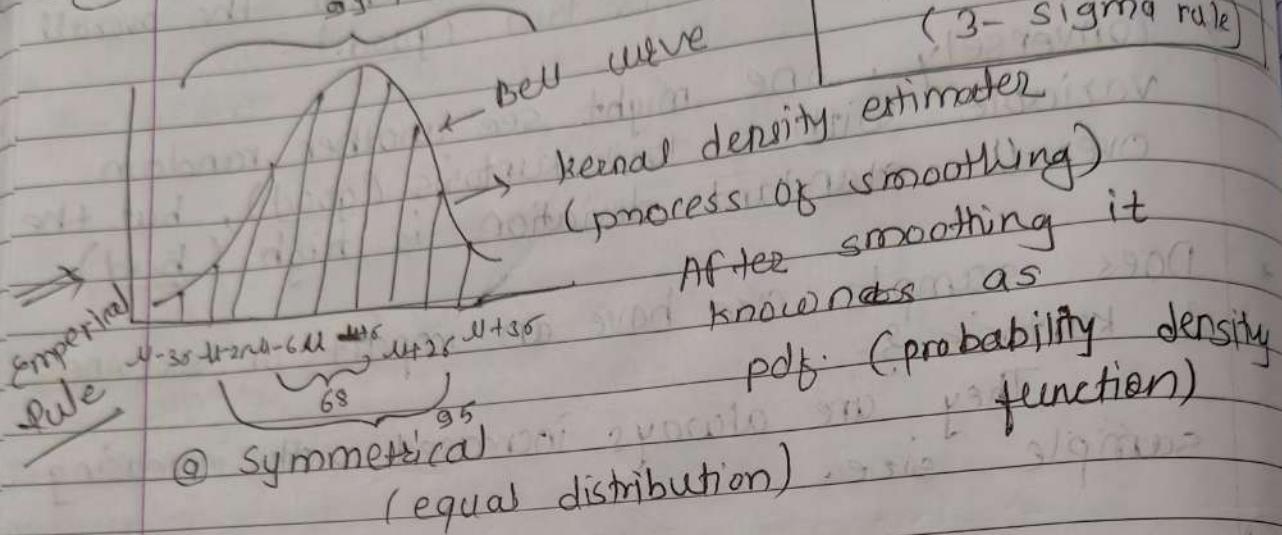
$$\frac{n(n+1)}{(n-1)(n-2)(n-3)} \leq \frac{(x_i - \bar{x})^4}{S.D.} - \frac{3(n-1)^2}{(n-2)(n-3)}$$

* Data Visualization Tools :-

How to construct a box plot for outliers :-



* Normal Distribution / Gaussian Distribution



In normal distribution follows 3 sigma rule.

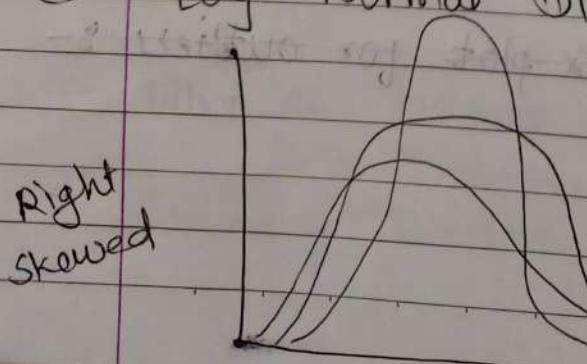
e.g. Height, weight

Iris Dataset \Rightarrow {sepal length, petal length}

* Central Limit Theorem :-

whether ~~whether~~ your distribution is normal distribution or ~~or~~ not normal distribution if we take the sample where $n > 30$ and if we take all the sample and when we ~~populat~~ or when we plot it we get a normal distribution

* Log Normal Distribution :-



if $x = \log$ normal distribution
then $y = \ln(x)$

Normal
 \log

$\log e$

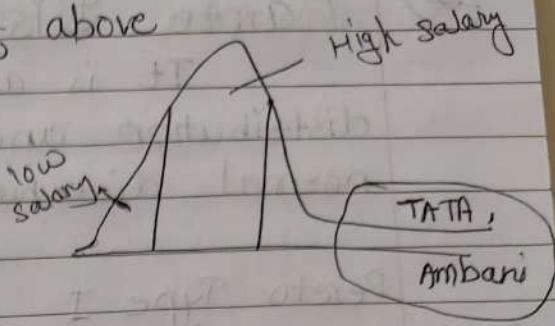
whatever value will get in y as it's coordinate with x it will be

Also, we can write instead of above

$$x = e^{a \ln(y)}$$



exponential.

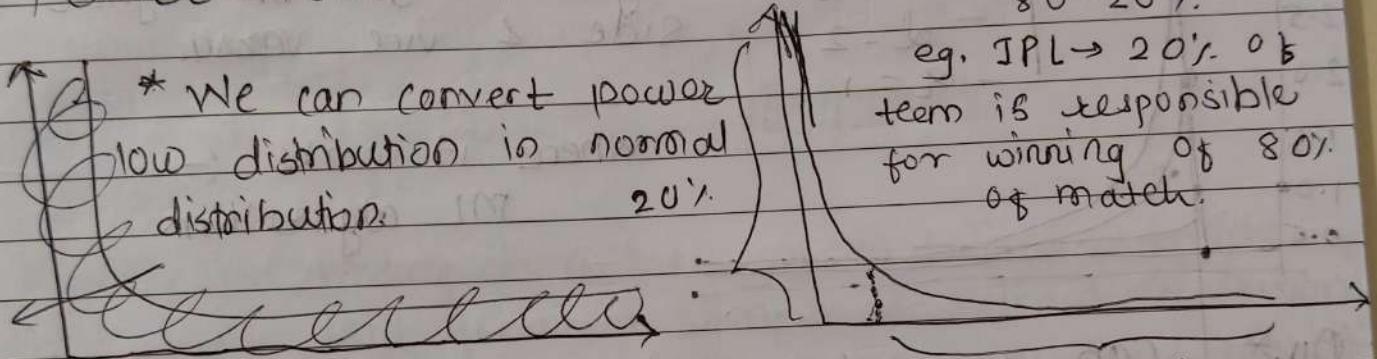


eg. a) wealth distribution

b) comments in

where we use s- ① machine learning

* Power law Distribution :-



eg. IPL → 20% of team is responsible for winning of 80% of match.

And the transform known as boxcox Transform.

eg. 80% of wealth is distributed with 20% of the total population.

80% of the total oil is distributed with 20% of nations.

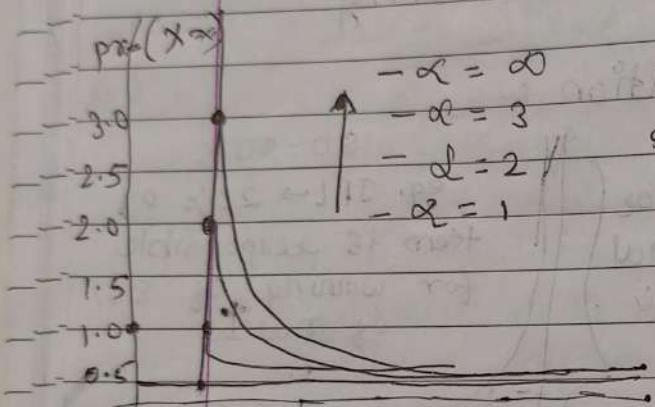
When we convert power law distribution & pareto distribution in normal distribution it's known as boxcox transform.

(*) Pareto Distribution :-

It is an example of power law distribution and this distribution are non normal distribution.

Pareto Type I

Probability density function.



As the element number goes higher it goes on higher side & vice versa.

where it works :-

ML algorithm

Qn.* Can we change above distribution in normal distribution.

Ans. Yes, we can transform in normal distribution

The process of transformation known as Box Cox transformation.

Eg. ① 20% of the products in Amazon is responsible of 80% of sales.

② 20% of defects solves the 80% of upcoming defects.

③ 20% of team is responsible in deliverable 80%.

Interview question :-

Relation betn log normal distribution & pareto distribution?

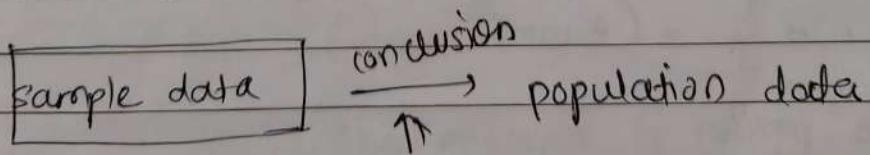
Hypothesis TESTING

Defⁿ Hypothesis testing is a form of statistical inference that uses data from a sample to draw conclusions about a population parameter or a population probability distribution.

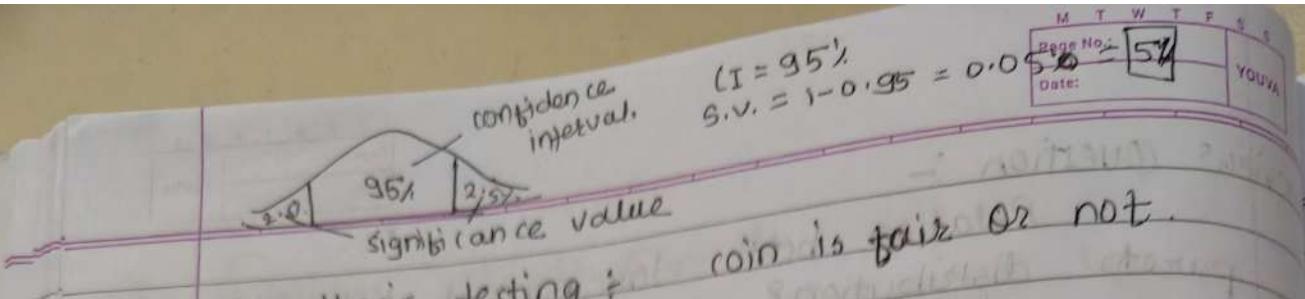
* What is hypothesis testing :-

Hypothesis testing is an art in statistics whereby an analyst tests an assumption regarding a population parameter. The methodology employed by the analyst depends on the nature of the data used and the reason for the analysis.

* Inferential statistics :-

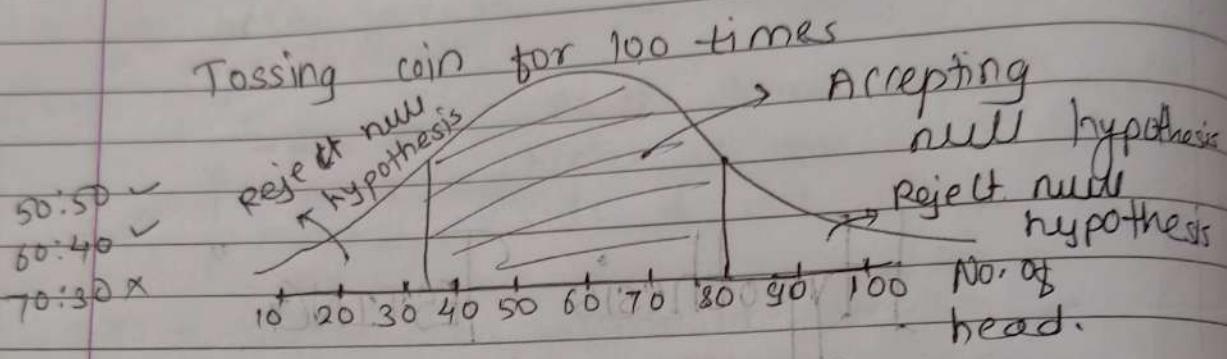


Hypothesis testing
(χ^2 test, t-test, Anova)



Hypothesis testing :- coin is fair or not.

- ① Null hypothesis (H_0) = coin is fair.
- ② Alternate hypothesis (H_1) = coin is not fair.
- ③ Experiment = where null hypothesis accepted or not accepted.



Domain
expect
decide
the

Confidence interval :- (percentage)

A confidence interval is a range of values used to estimate an unknown statistical parameter, such as a population mean. (Decision boundary)

* 2 Tail Test :-

A two-tailed test is a hypothesis test where the alternative hypothesis claims that the population parameter is different from the value specified in the null hypothesis, meaning it could be either greater than or less than.

* Distribution function :-

- ① Probability distribution function
- ② Probability Density function
- ③ Probability mass function
- ④ Cumulative distribution function

We understand distribution of Probability distribution function data.

Probability density function

\uparrow distribution

eg continuous Random variable

Probability mass function

\uparrow distribution

eg Discrete random variable (whole number)

$$\text{a) Variables :- } x + 5 = 7$$

$$x = 7 - 5 = 2$$

$$3y + x = 8$$

$$y = \frac{8 - x}{3} = 2$$

In linear algebra we used to find x & y , etc. So, these are nothing but the variables.

b) Random variable :-

Is a process of mapping the output of a random process or experiments to a number.

eg ① Tossing a coin (experiment) \Rightarrow outcome of experiment.

The values
in variable
changing

$$X = \begin{cases} 0 & \text{if head} \\ 1 & \text{if tail} \end{cases}$$

② Rolling a die (experiment) = $\{0, 1, 2, 3, 4, 5, 6\}$

$$Y = \begin{cases} \text{sum of the outcome of rolling a die 5 times} \end{cases}$$

$$[\text{probability } Y > 15] = \underline{\text{to find}}$$

Interview question

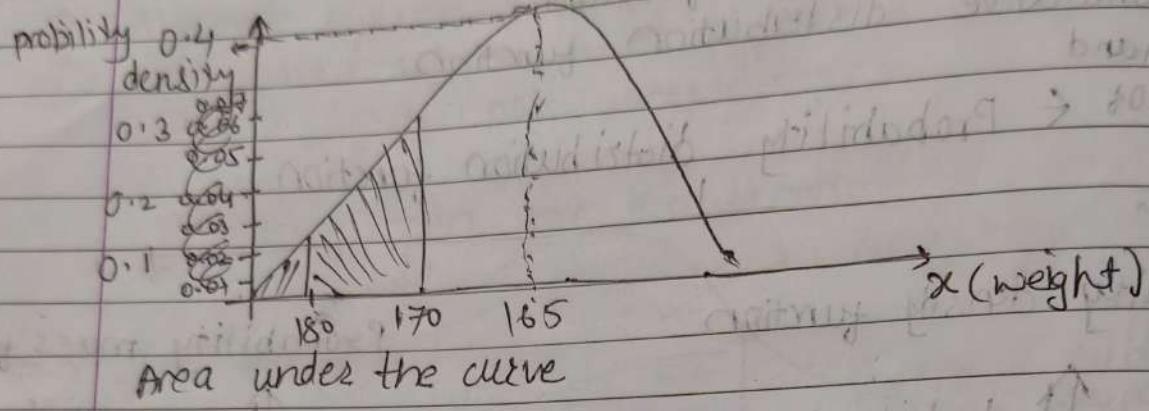
Difference betn Pdt & Cdf

Date:

YOUVA

(*) Probability Density function (PdF)

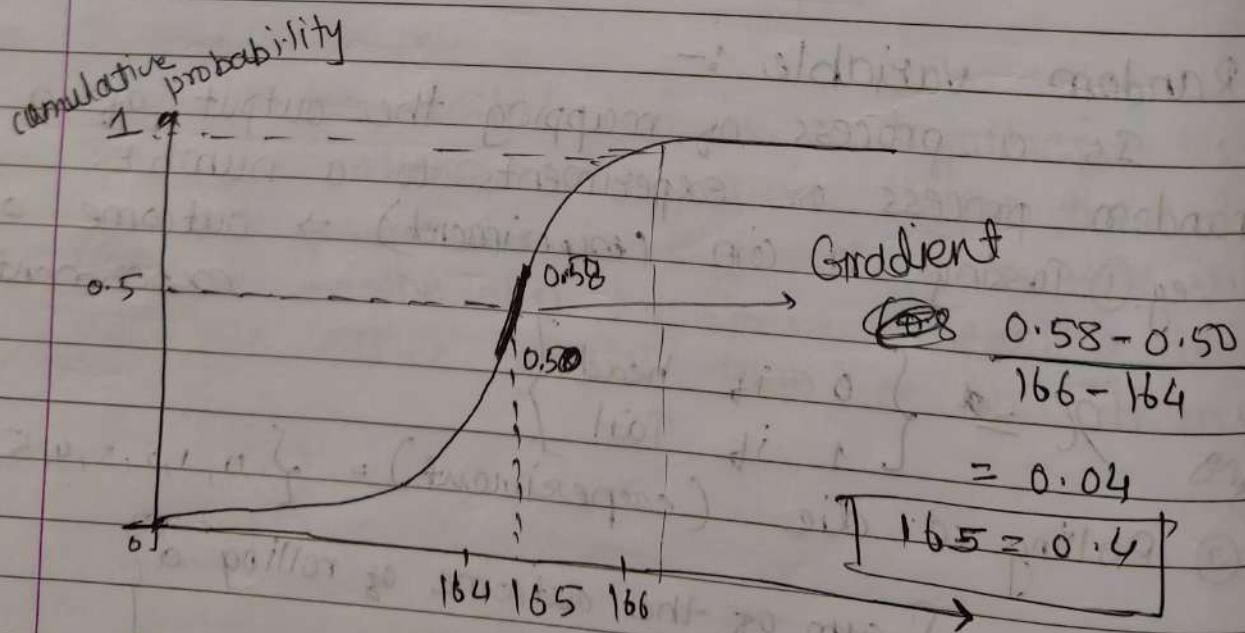
(@) Distribution of continuous Random variable



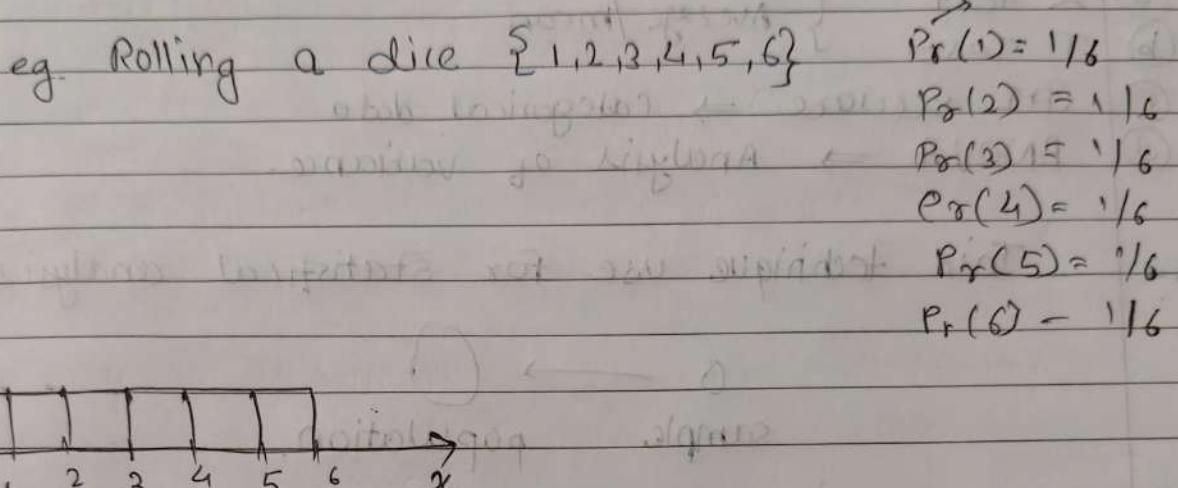
probability Density = Gradient Percent of cumulative curve.

Gradient =
moment in the
x axis
moment in the
y axis

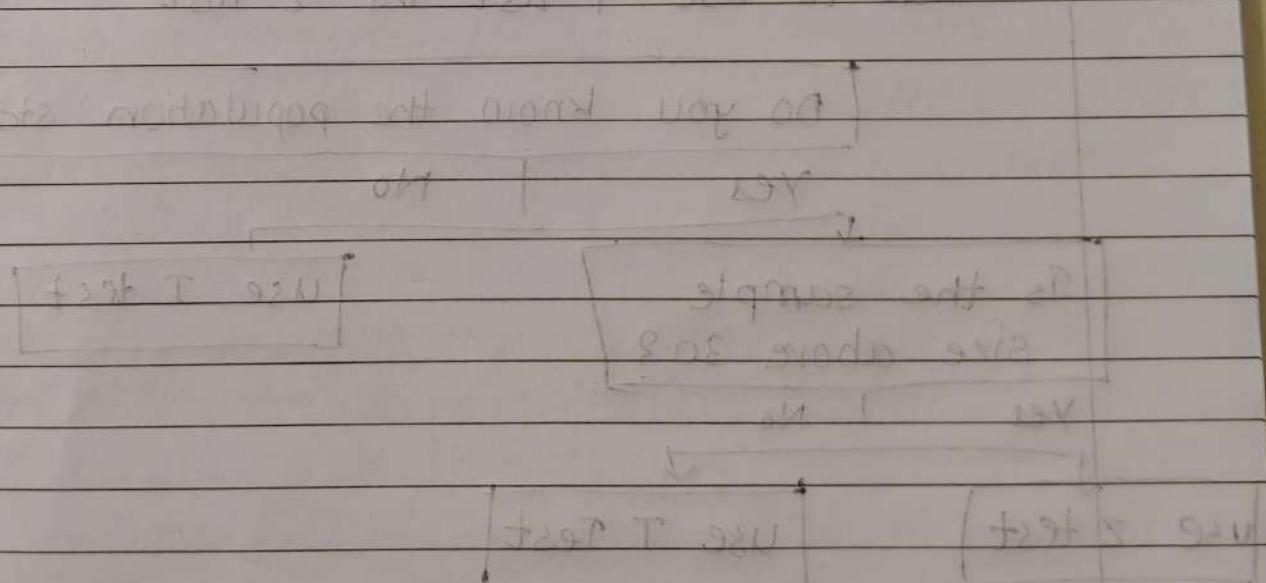
* Cumulative Density function :- (CDF)



* Probability Mass function of Discrete Random Variable.



* Difference b/w PMF & CDF

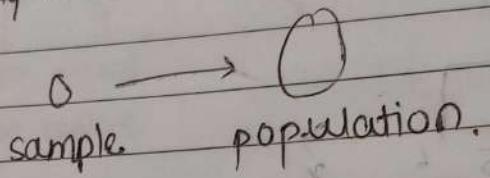


in a n histogram fo frequency approximatly
probabilitiy distribution is often in the form of
binomial id of norm with correlated variable A
+ Uniquelikely to go equal with binomial in
case of 2.01 ad. etc. because out bound

* Hypothesis Testing & Statistical Analysis

- (a) Z test } Average / mean
- (b) T Test
- (c) Chi square \Rightarrow categorical data
- (d) ANOVA \Rightarrow Analysis of variance.

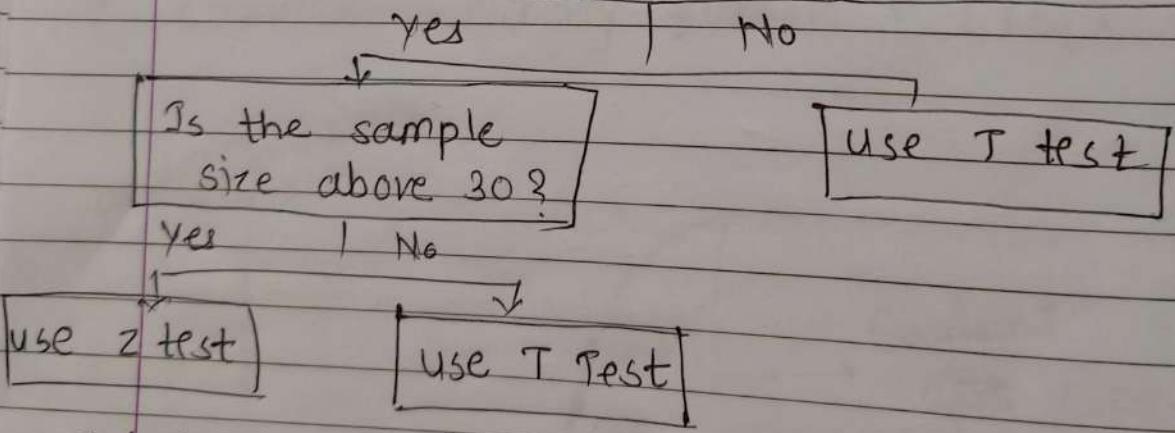
This technique use for statistical analysis.



- (a) Z test :- Z Table

When to use T test Vs Z test

Do you know the population std. σ ?



* Z-test

The average height of all residents in a city is 168 cm with a population std $\sigma = 3.9$. A doctor believes the mean to be different. He measured the height of 36 individuals & found the average to be 169.5 cm.