

Statistics is the science of collecting, organizing and Analyzing data.

Data :- facts or piece of information.

Eg:- Height of students in a Class

140 cm, 145 cm, 146 cm, 180 cm, 172 cm.

Eg:- Gender of person visiting a doctor (8am - 4pm)  
(M, M, F, T, F, M)

Statistics is of Two types  $\rightarrow$  Descriptive Stats.  
 $\rightarrow$  Inferential Stats

Descriptive Stats

$\hookrightarrow$  It consist of organizing and summarizing of data.

Inferential Stats

$\hookrightarrow$  It consist of using data you have measure to form conclusion.

D-S  $\rightarrow$  measure of central tendency (Mean, Median

$\rightarrow$  measure of dispersion

$\hookrightarrow$  Variance

$\hookrightarrow$  Standard Deviation

$\rightarrow$  Different types of Distribution of Data

$\rightarrow$  Histogram, pdf, pmf.

I.S  $\rightarrow$  Diff of Type of Test

- $\hookrightarrow$  Z Test
  - $\hookrightarrow$  t-test
  - $\hookrightarrow$  Chi square test
  - $\hookrightarrow$  ANOVA test.
- } Hypothesis testing.  
H<sub>0</sub>, H<sub>1</sub>, p-value, Significance value.

Q: Let consider there are 20 classroom in a college and let say you have collected the ages of student in one class.

$\rightarrow$  Ages of class : {19, 21, 18, 34, 22, 21, 25, 20}

Descriptive Stats Questions.

$\hookrightarrow$  What is the common age in your stats class.

Ans  $\rightarrow$  Mean (Ages)

$\hookrightarrow$  which is a<sup>part</sup> central tendency.

Inferential Stats Question.

$\hookrightarrow$  Are the ages of student in the classroom

similar to what you expect to

the ages of the student in the University.

We have concluded  
that the ages are  
similar or not  
Sample

In I.S, we saw that sample & population is going to be performed.

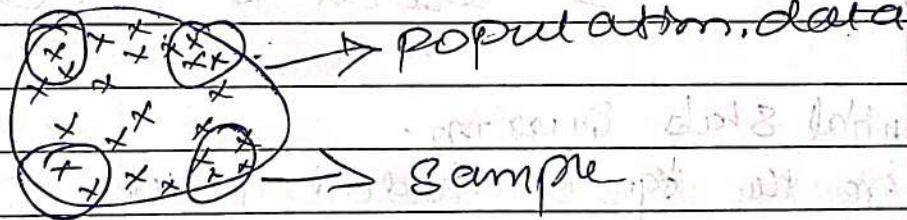
## Test 2 Population & Sample, And Sampling Techniques

Population :- The group that we are interested in studying.

Sample :- It is a subset of a population.

Eg. :- EXIT POLL → Assumption made on the basis of sample data, that what are the chances of winning the state election for ~~the~~ party which

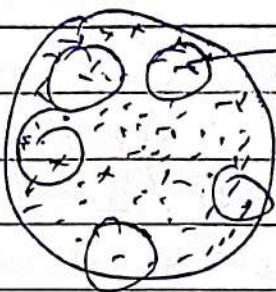
STATE A (2014) Age  $>$  18 year



So, On the ~~basis~~ By reading out to the sample data and getting the insight through it, will make an assumption ~~on~~ on the basis of maximum vote winning party has the chance of winning the election.

But the EXIT POLL, it is not always true. This is just an assumption.

STATE B, to average height of the people.

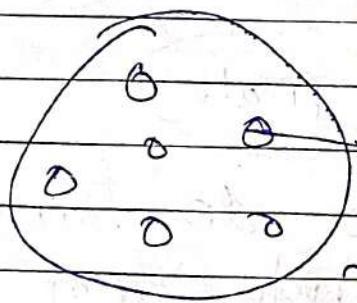


Note → As we choose sample from the population, But what are condition/math or process of taking the samples from the population. what is the criteria behind it, so to know that we should know Sampling Techniques

Sampling Techniques ⇒ The main motive of this ~~is~~ is, that when we choose sample. So on what base we are choosing it.

The goal of Sampling is to create a sample that is representative of the entire population.

Population is denoted by (N)



Sample:

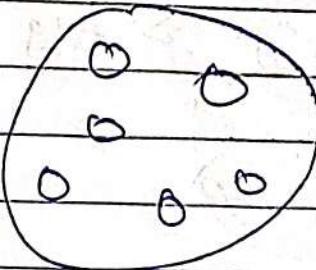
This entire Sample represents population.

Sample is denoted by (n)

Let's, find out how many types of Sampling techniques are there.

## Types of Sampling.

### ① Simple Random Sampling.



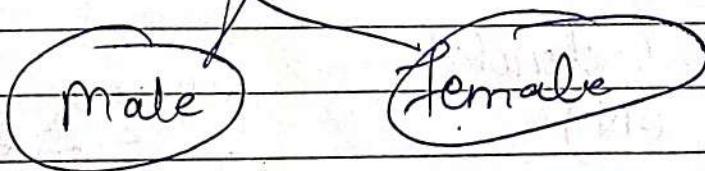
Randomly we are taking any samples without any condition applying to it.

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:

Stratified means layering, we also define it as ~~state~~ strata. And these layering are non-overlapping groups.

Eg:- Population



If I wanna do Survey on Male, then I need info from male gender, so I'll reachout to Male.

Same goes for female.

Also here, you multiple groups, but make sure that these groups shouldn't

### ③ systematic Sampling :-

Assume, we are doing a survey outside a Hall every 4th person, I'm asking for to fill the needful for the Survey.

Whenever you decide for sampler, you decide select a team "4" (1st) and whenever I see the nth person, I'll approach to that person for the Survey.

### ④ Convenience Sampling:- If we want to do Sampling, expertise or there is term called "Voluntary Response Sampling"

This term is very important with Convenience Sampling.

Eg: I want to do a survey on Data Science.

With Approach  
to them, those

who have knowledge in DS,  
expert, student in DS.

i.e. Can't approach Doctor. Right!

Data Science

Exercise!  
→ EXIT POLL (Random)

→ Disease Information (Convenience Sampling)

→ Household expenses  
(Strategic Sampling)  
Main (Women)  
Target

TUT 8

## What are Variable & Its Types.

Defn: Variable is a property that can take on multiple / many values.

Eg:- age = 12 height = 172 cm  
13 172.5 cm  
40 180 cm

Weight = 72 kg, 72.5 kg, 73 kg.

Variable  $\rightarrow$  Singular mode.

whenever we talk about variable we talk  
with respect to only with a single value.

Ages = [12, 24, 48, 60, 100]

↳ plural mode.

## Types of Variable

## Quantitative Variable

Qualitative  
Variable  
(Categorical)

Discrete Variable

Continuous  
Variables

Based on the properties  
This Variable is further

$\Rightarrow$  The value where Discrete variable will be assigned will always be a whole Number

Eg:- Tuftnile

## Values.

Eg & Wedge

E.g.: Gender  Male  
 Female

It is divided on some properties.

## Types of Flower

~~rose~~

L. L. C.

## Cœurs.

② Median

③ Mode

Defn: Central Tendency refers to the measure used to determine the "center" of the distribution of data.

And (mean, Median, Mode) are the techniques to find out the center of distribution of data.

Eg:

{1, 2, 2, 3, 4, 5} lets find out Mean (Average)

Usually when we talk about (mean, median, mode) we mostly focus on two types of Data  $\rightarrow$  population Data ( $N$ )  $\rightarrow$  Population Mean  $\rightarrow$  Sample Data. ( $n$ )  $\bar{x}$   $\rightarrow$  Sample Mean

Population Data ( $N$ )

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

Sample Data ( $n$ )

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

$x = \{1, 2, 2, 3, 4, 5\}$  Data  $\rightarrow$  Population

let find the average (Mean)

## ② Median

$$X = \{1, 2, 2, 3, 4, 5\} \quad \bar{x} = \frac{1+2+2+3+4+5}{6}$$

$$\bar{M} = 2.83$$

found.

## ② Median

$$X = \{1, 2, 2, 3, 4, 5, 100\}$$

This number is called outliers, which is out of scope.

$$\bar{M} = 2.83$$

$$\bar{M} = 2.83 \quad (\text{Previous Mean})$$

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

$$\bar{M} = 16.71$$

There is a huge diff in both of this  $\bar{M}$  (Mean).  $16.71$  is out of scope just because  $100$  got added, which is a case of outliers, so, to get rid of these outliers we'll use median.

Median says there are 3 steps

① Sort all the numbers

② find the central element

$$\begin{array}{l} \xrightarrow{\text{Odd length}} \{1, 2, 2, 3, 4, 5, 100\} \\ \xrightarrow{\text{Even length}} \end{array}$$

Odd length  $\rightarrow$  Cause there are 7 elements  $\rightarrow$  median = 3

Even length / observation

$$\{1, 2, 2, 3, 4, 5, 100, 101\}$$

$$\frac{3+4}{2} = 3.5$$

⑧ Mode :- Most Frequent Elements.

If check the frequency of the data element of String data / Numerical data.

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

most frequent Element.

We use Mode, when we have a Categorical feature.

Eg:-

	Age	Weight	Gender
Missing value	24	-	M
and	25	-	F
get	-	78	
replace	-	80	
by	26	82	
the	27	-	
mean	-	-	

Suppose Here Male is repeated value, then the missing value will get replace by M - male.

(Weight, 9)  
in case

there is

outlier with

use

Median  
Based

on

Even or

odd observation

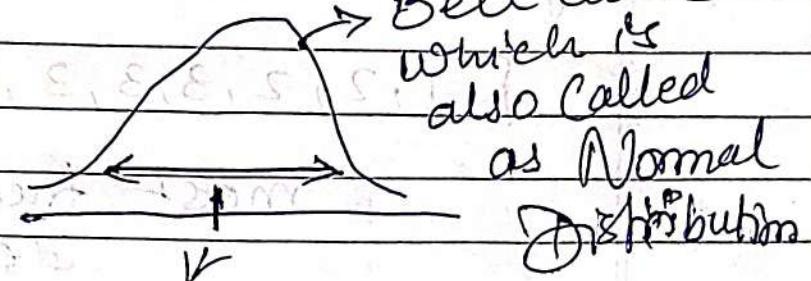
# T.S Measure of Dispersion - Variance

(a)  $\downarrow$

**Spread**

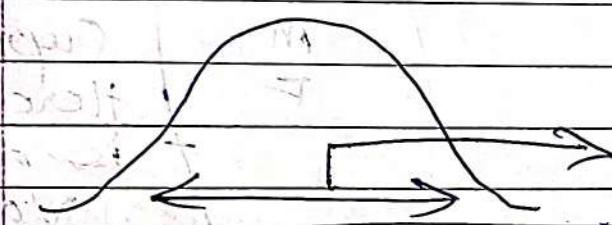
Standard Deviation.

- ① Variance
- ② Standard deviation



When ever we speak about this Central element, it can be either our Mean, median, mode.

↳ And it specifies measure of central tendency



This Spread, which is also called as Dispersion

And this spread is determine using two important things → Variance  
ii. → Standard Deviation

population ( $N$ )

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N (x_i)$$

Sample ( $n$ )

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

Population Variance

σ² → Sigma

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

Sample Variance

s² → "s" square

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

Degree of freedom

Bessel Correction

$$\text{Q: } X = \{1, 2, 2, 3, 4, 5\}$$

$x_i$	$\bar{x}$	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
2	2.83	-0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71
$\bar{x} = 2.83$			10.84

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

$$= \frac{10.84}{6-1}$$

$$= \underline{10.84}$$

$\sqrt{s^2}$  Spread of the

Sample Variance  $\rightarrow$  Data Dispersion  
 $\sqrt{s^2} = \sqrt{10.84} \rightarrow$  we can say)

Standard Deviation formula

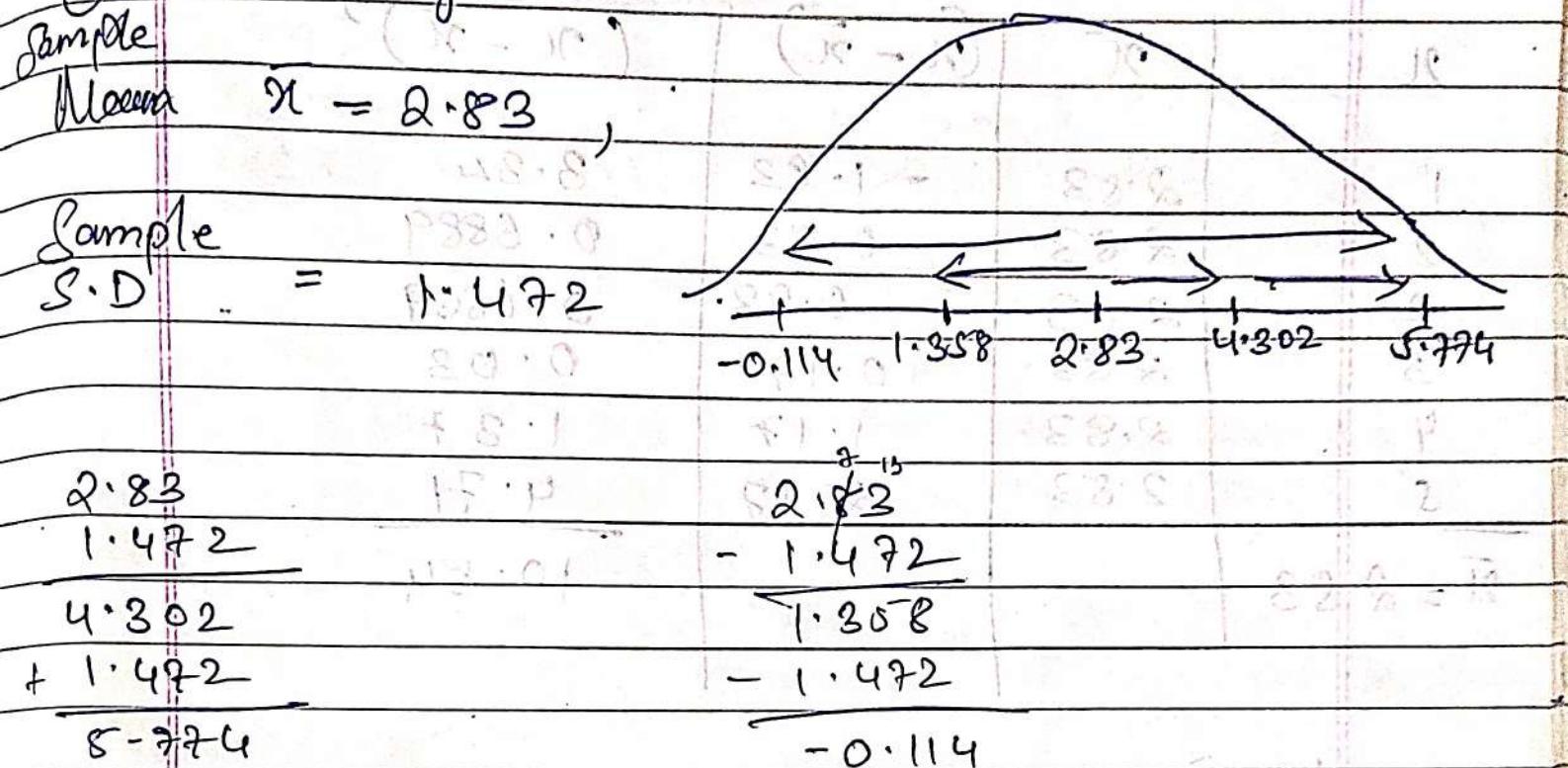
$$\hookrightarrow s.d. = \sqrt{\text{Variance}}$$

$$= \sqrt{10.84}$$

Sample

$$= \boxed{3.29} \rightarrow \text{Standard Deviation}$$

# Dataset Insights



Here you got the idea of Standard Deviation  
But still you don't know what is the work  
of Role of Variance.

Variance decides the entire spread or dispersion.

lets say  $\rightarrow$  Variance - Big no. Variance - Small no.  
S.D will also be Big S.D will also be small no.

Spread / Dispersion  
will be more

① Note  $\rightarrow$  if Variance is a big number  
spread is also increasing.

② if Variance is small number  
spread is also decreasing.

Spredness  
will be less,  
But height  
will be more.

## T.6 Percentile and Quartiles

Date \_\_\_\_\_

Dif betw Percentage & Percentile.  
Percentage = {1, 2, 3, 4, 5}

% of all numbers that are even?

$$\% \text{ of even} = \frac{2}{5} \times 100 = 40\%$$

So this was percentage

Now, Percentile:- A percentile is a value below which a certain percentage of observations lie.

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

Dataset :- 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

What is the percentile ranking of 10?

Step 1: Note how many numbers are there in a dataset  $n = 20$

$$\text{Percentile Rank of 10} = \frac{\text{total no. of value below } 10}{\text{total no. of observations}} \times 100$$

Now, what does this  
80 percentile  
indicates

The value 10 is 80% bigger value than entire distribution value.

Now for 12, find the percentile value for 12.

$$\text{P} = 0.95 \quad x = 12 \quad n = 20$$

Percentile Rank of 20 = total no. of Obs below  $x \times 100$

$$\frac{19}{20} \times 100$$

Value is 95 percentile

The value 12 is 95% bigger value than entire distribution.

Different Question

What value exists at percentile rank of 25. What value exist at 25 percentile

formula :-

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

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

$$= 0.25 \times 21$$

$$= 5.25$$

5.25 is not the value it is the index from the given dataset for

Glen

~~given dataset~~ →  $2, 2, 3, 4, 5, 5, -5, 6, 7, 8, \dots$

which is  $\frac{1}{5}$  but  $5,5$

$$\rightarrow \frac{5+5}{2} = \frac{10}{2} = 5$$

So, 5 is the value, which lies at 25 percentile.

## Quartile

## Box Plot

$\rightarrow$  1st Quartile (Q1) - 25%

Median  $\rightarrow$  2nd Quartile ( $Q_2$ ) - 50%

→ 3rd Quartile (Q3) - 75%

## Inter Quartile Range

$$\Rightarrow \text{Lower Bound (LB)} = Q1 - 1.5 \text{ (IQR)}$$

$$\Rightarrow \text{Upper Bound (UB)} = Q3 + 1.5 \text{ (IQR)}$$

Quartile 18

To find outliers, we specifically use Box plot.

To know how Box plot is created, we have to use Quartile

IQR (Inter Quartile Range)  $\Rightarrow$  And from this IQR we find the outliers.

## T.F How to construct a Box Plot for Outliers.

But - Before that we'll solve a problem. (Chaitanya Sir)

Dataset

$$\Rightarrow \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27\}$$

First will find Q1 - 25% percentile

formula  $\rightarrow$

$$Q1 = \frac{\text{percentile}}{100} \times (n + 1)$$

$$\frac{25}{100} \times (19 + 1)$$

$$\frac{25}{100}$$

Index = 5

50

Index = 10

Index = 5th Index, Value = 3

Index = 15

Now, Q2, - 50% percentile

$$\frac{50}{100} \times (n + 1)$$

Q2 = 25th Percentile

$$\frac{50}{100} \times 25$$

Q2 = 10th Index, Value =

Now  $Q_3$  - 75% percentile

$$Q_3 = \frac{\text{percentile}}{100} \times (n+1)$$

$$\therefore Q_3 = \frac{75}{100} \times (20)$$

$\Rightarrow$  15th index, value = 7

Now IQR (Inter Quartile Range)

$$\begin{aligned} IQR &= Q_3 - Q_1 \\ &= 7 - 3 \end{aligned}$$

$$IQR = 4 \rightarrow IQR$$

Now LB (Lower Bound)

$$\begin{aligned} LB &= Q_1 - 1.5(IQR) \\ &= 3 - 1.5(4) \\ &= 3 - 6 \\ &= -3 \end{aligned}$$

And UB (Upper Bound)

$$\begin{aligned} UB &= Q_3 + 1.5(IQR) \\ &= 7 + 1.5(4) \\ &= 7 + 6 \\ &= 13 \end{aligned}$$

$$UB = 13$$

## → Five Nos Summary And Box Plot.

① Minimum

② First Quartile (25%) Q1

③ Median

④ Third Quartile (75%) Q3

⑤ Maximum

All these Components are used to Create Box plot.

And the main aim of this Box Plot is that it helps to detect the Outliers

~~So, we can remove the outliers.~~ Data should be removing the outliers. In sorted.

Dataset  $\rightarrow \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27\}$

→ So, to detect the outlier, this technique is most favourable in Machine Learning.

Find,

Lower fence  $\longleftrightarrow$  Higher fence

Now, lets 27 is a outlier, which is extremely different from all the distribution value.

If lower fence is 5, or if can be any value . so below the lower fence value .

whatever number lies will be Outliers , And for higher fence = 30 or any value , whatever value above higher fence value lies

But we cannot just confirm that its an outlier , Because there would be lot of data .

Even there could be

$$\text{Lower fence} = Q1 - 1.5 \times (\text{IQR})$$

$$\text{Higher fence} = Q3 + 1.5 \times (\text{IQR})$$

$$Q1 = \frac{\text{percentile}(25)}{100} \times (n+1)$$

$$Q3 = \frac{\text{percentile}(75)}{100} \times (n+1)$$

LB —  $\circlearrowleft 3 \circlearrowright \leftrightarrow 13 \circlearrowleft \text{UB}$

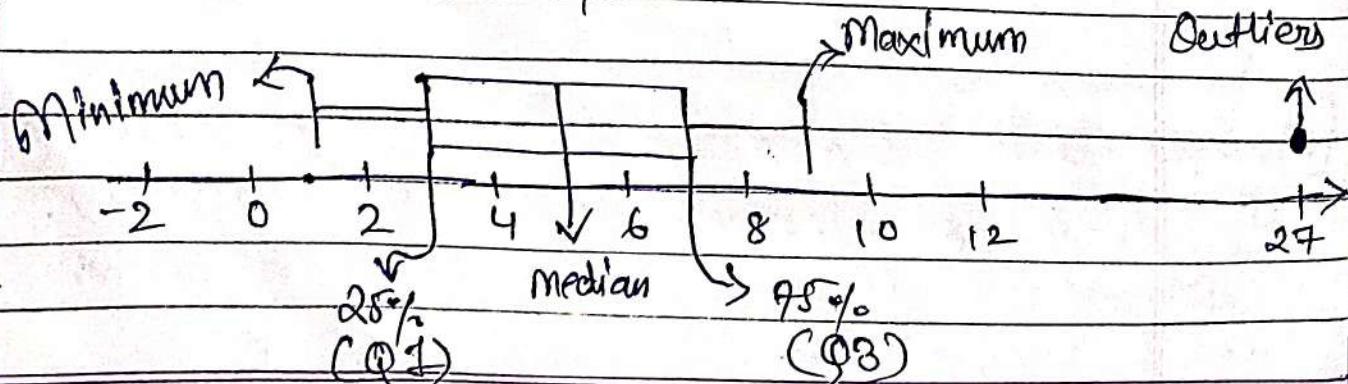
This is the range from  $-3$  to  $13$

So,  $27$  is definitely an outlier. So we can remove it or either we can plot it with a Box plot.

So, let's move to the 5 components, as per the dataset

- ① Minimum value = 1
- ② first Quartile (25%)  $Q1 = 3$
- ③ Median = 5 (If we calculate, center of elements)
- ④ Third Quartile (75%)  $Q3 = 7$
- ⑤ Maximum = 9 (we won't consider outliers)

Now, we'll Create Box plot.



apply Z-score, And this process is called Standardization.

Z-score is It tells how much standard deviation it is away from the mean.

formula :-

$$Z\text{-score} = \frac{x_i - \bar{x}}{s} \rightarrow \text{standard deviation.}$$

X	$\bar{x}$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	working.
1	2.83	-1.83	3.35	1 + 2 + 2 + 3 + 4 + 5.
2	2.83	-0.83	0.68	
2	2.83	-0.83	0.68	
3	2.83	+0.83	0.17	0.03
4	2.83	-2.83	1.17	1.36
5	2.83	-3.83	2.17	4.70
$\bar{x} = 2.83$				10.80

$$\text{Mean } (\bar{x}) = 2.83$$

Finding Variance

$\hookrightarrow$  Variance

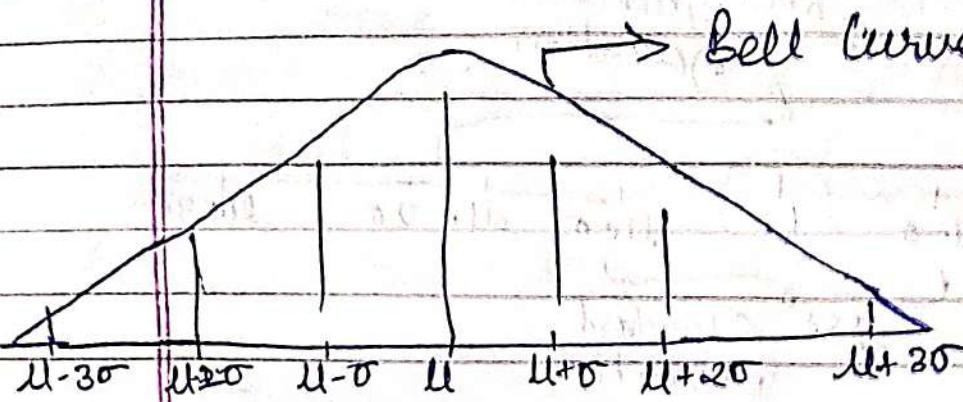
$\sqrt{\text{formula}}$

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

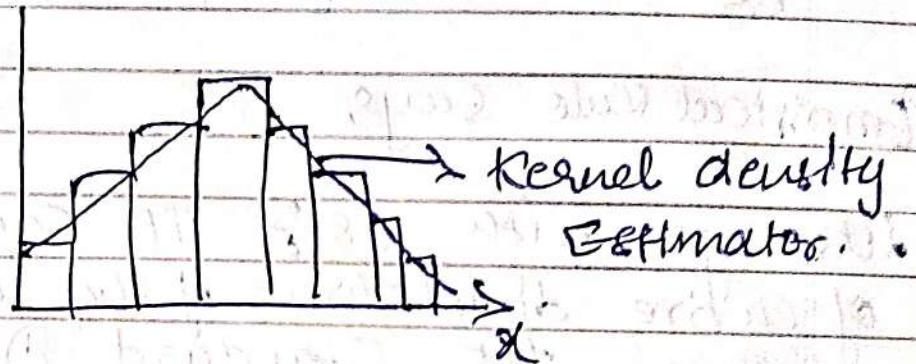
$$\sigma^2 = \frac{10.80}{6}$$

T.9

## Normal Distribution / Gaussian Distribution



To know it ~~as~~, in simple term that how a Bell Curve is made will take ~~eg~~ of histogram.   
~~eg~~: let's take an histogram,

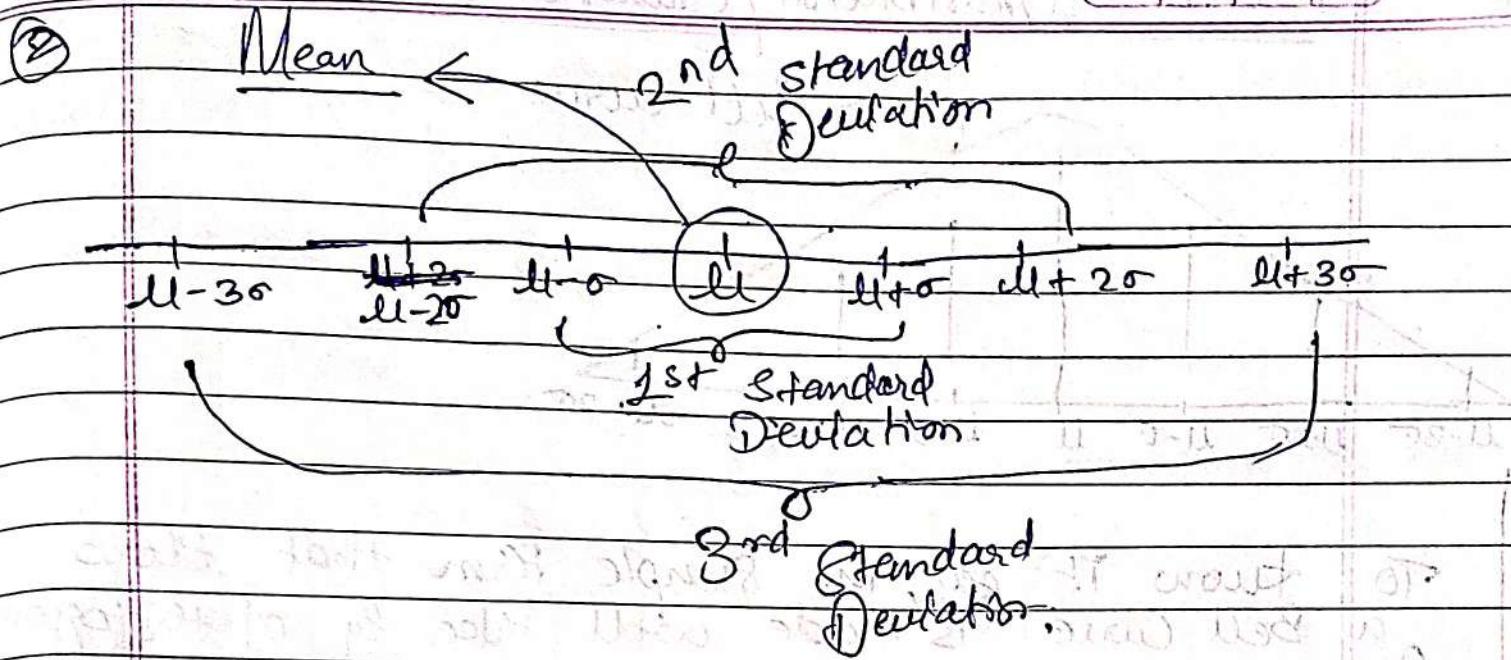


the Bell curve is apply after smoothing process, And this smoothing is called kernel density estimator.

And this kernel density, converts this histogram into pdf (probability density function).

And after doing this process, we get this Bell curve.

→ ① Symmetrical - [Left area and Right area of the curve are Both Equal.]  
for Eg:- we have data set: {1, 2, 3, ... }  
\* Symmetrical means 50% of number is on the left side of mean and remaining 50% is on Right side of the Mean.



**68 - 95 - 99.7 [Empirical Rule]**

Empirical Rule says,

lets go with 68% It says that 68% of entire distribution of the data lies within 1st Standard Deviation.

let assume, we have a dataset of 100 elements. so, 68 data element will lie within 1st standard Deviation as per Empirical Rule.

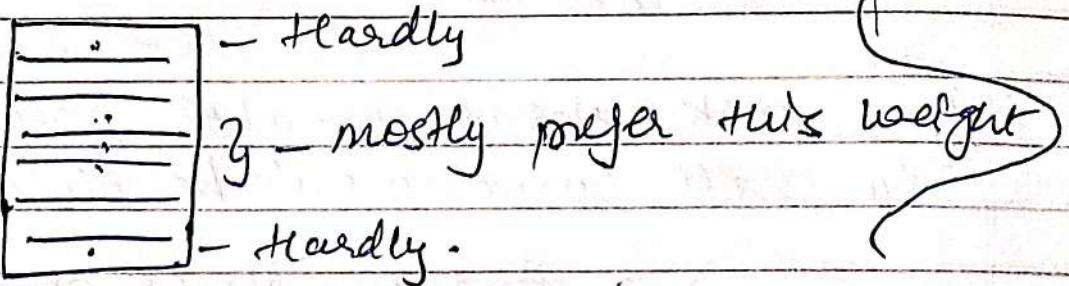
Now 95, 95% of the data of this distribution following falling within 2nd Standard Deviation.

And 3rd Standard Deviation of Data of this distribution covers 99.7% of data. only if it follows Gaussian Distribution or Normal Distribution.

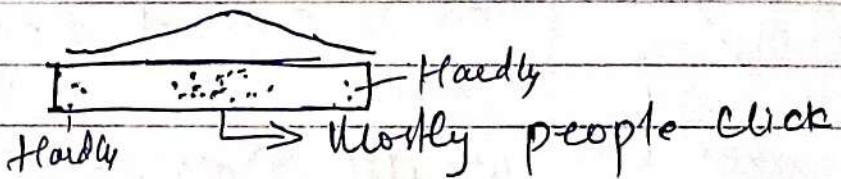
Note: If the distribution is symmetric and bell-shaped, then definitely it follows the 3-sigma rule of Empirical which are [68 - 95 - 99.7]

Let's see some examples, like what kind of dataset does normal or Gaussian Distribution follows:

Eg: Gym weight plates



Space bar



Note:

Whenever it is Gaussian Distribution or Normal Distribution then definitely it follows the 3-sigma rule of Empirical rule [68 - 95 - 99.7] and this is an important property of Gaussian or Normal Distribution. and with the help of this you can solve many problem as

- ① you can find Outliers (By the 3rd standard deviation, like whatever data lies outside the 3rd standard Deviation will be Count as Outliers.)

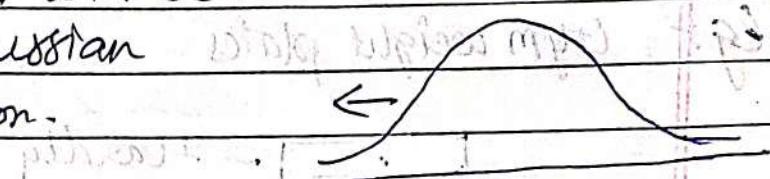
## T10 Central Limit theorem -

Question  
Date \_\_\_\_\_  
Page \_\_\_\_\_

weight Data ( $x$ ) =  $65, 72, 83, 55, 64, 67, \dots$   
Variable ← weight of a population.

So, the Bell curve will be

Normal / Gaussian Distribution.



If I talk about wealth Data, then the bell curve would be like,

→ And this Right Skewed Distribution is called Log Normal Distribution

wealth Data.

Sample size → If our population is normal / gaussian Distribution or lognormal Distribution and if we apply some Sample size (cond<sup>n</sup>) randomly to the distribution data then we'll always get a Normal / Gaussian Distribution.

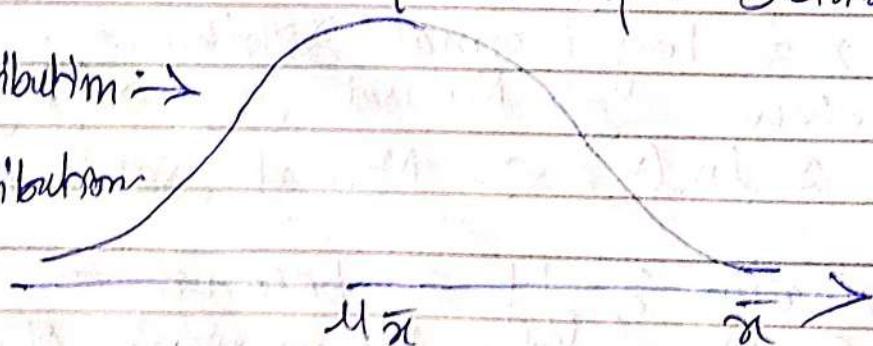
Sample size → Choosing random set of samples from the population Data and with the set of sample Mean plotting the graph we'll always get a Normal / Gaussian Distribution.

So, Central Limit Theorem tell that the distribution is Gaussian or Normal Distribution or log Normal Distribution or whatever type of Distribution it may be, if we select a sample size which was  $n > 30$  and we take Mean Sample Mean of each set of Sample and plot it on the graph then it will always form Normal / Gaussian Distribution.

And this Central Limit theorem says whether your distribution is Gaussian or not, it doesn't matter. ~~cause if it takes~~ cause if it takes Sample Size condition of  $n > 30$ , then it will always give Normal / Gaussian Distribution.

→ Sample Distribution.

Normal Distribution →  
or  
Gaussian Distribution



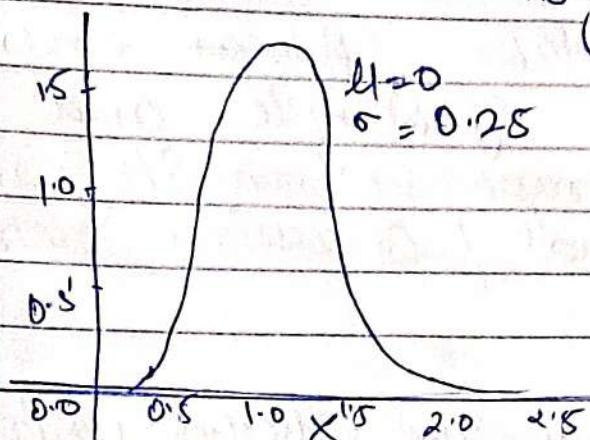
See Kolmogorov Central Limit Theorem.

## T-II Log Normal Distribution.

Page No.

Date

- we'll study  $\rightarrow$  what is log Normal Distribution
- $\rightarrow$  what its properties
- $\rightarrow$  what its mathematical Eqn
- $\rightarrow$  Examples of Data which follow Log Normal Distribution-



$\rightarrow$  This is Called Right Skewed and this is an important property of log normal Distribution.

It seems to be like Gaussian Distribution type, But the Right side is elongated.

Let see this Distribution in Mathematical Way,

$x \approx$  Log Normal Distribution.  
then  $y \approx \ln(x) \Rightarrow$  Normal Distribution.

Explanation: If  $x$  belongs to Log Normal Distribution. (Take the example of above distribution graph) and we consider that it is a Log Normal Distribution.

So, to verify it we take a new random variable ( $y$ ) in which we have applied log to the base "e"  $\Rightarrow \ln e$  in every  $x$ 's variable, so whatever output or if I plot it, I'll get Normal Distribution.

And if this property get satisfied, then we can say that  $X$  belongs to Log Normal Distribution.

We can write that Algorithm Inversely also:-



$$X \sim \exp(y)$$

$$y \sim \ln(x) - \text{Previous Algorithm}$$

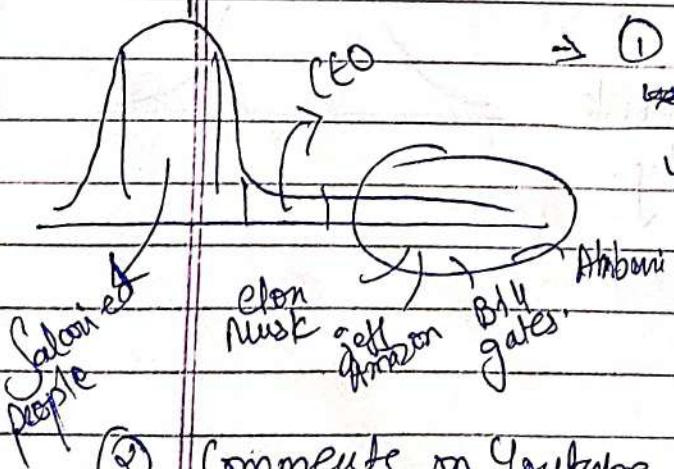
$X$  value ho Jayega

These Both are same

And If I inverse  $\ln$  then Obviously it will become exponential of  $y$   $\exp(y)$  and  $X$  will be come Value.

And these both will follow Normal Distribution  
if my  $X$  is Log Normal Distribution.

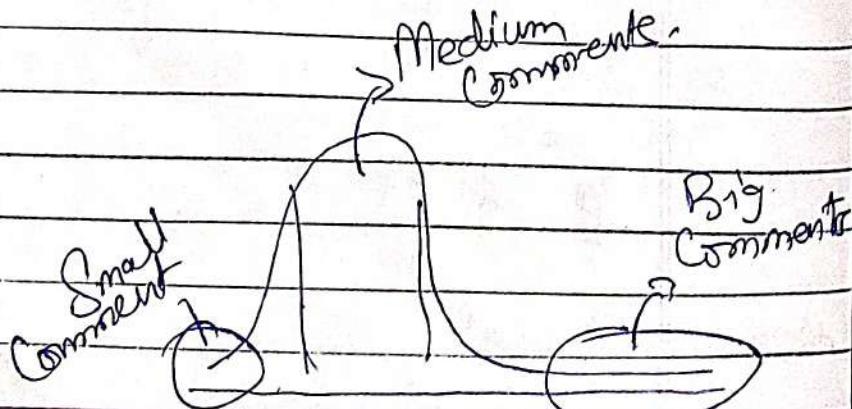
Examples:- let see the examples.



① wealth Distribution

Very less people have more wealth.

② Comments on YouTube



Now, Where we use this log Normal Distr

① Machine Learning Application  
When ever we in Machine Learning  
Algorithm we use dataset

If I work with simple linear regression

→ Normal Distribution, It will  
work effectively.

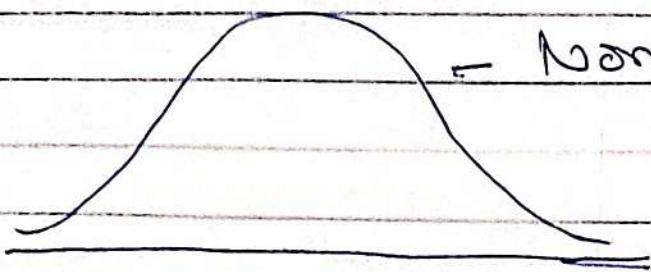
Right skewed, we'll apply loge  
→ If its X's variable  
and it will convert  
it into a Normal  
Distribution

And then it will work  
effectively or properly.

# T.12 Power law Distribution

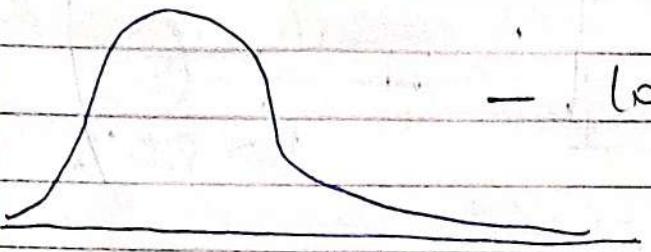
Date \_\_\_\_\_

①



- Normal Distribution

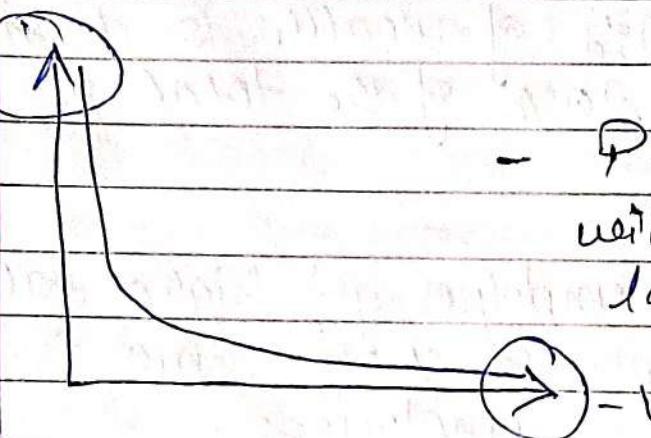
②



- log Normal Distribution  
(Right Skewed)

③

Variable 2



- Power Law Distribution,  
we'll see its examples  
later



- Variable 1

If you want to know about Power law Distribution  
Quickly then always remember this Rule

12

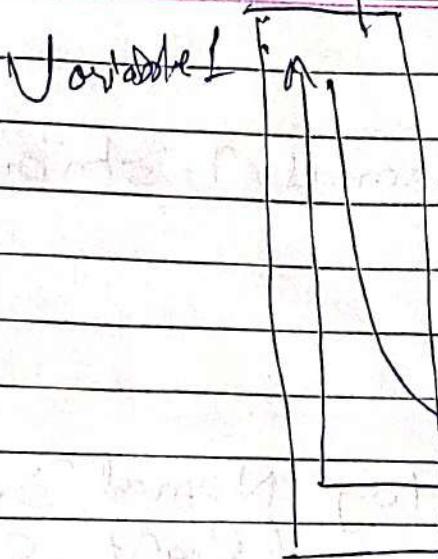
80 - 20% Rule

Both of these Variable forms a relationship  
of 80-20% Rule.

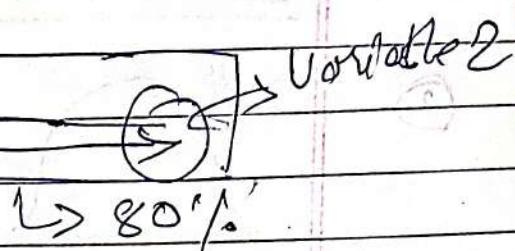
Eg: IPL  $\rightarrow$  20% of Team is responsible  
for winning 80% of match

Explanation: If you take any match of IPL, you  
20% of Team hangs Matlab 11 mai se  
2 ya 3 log itna Achha Kuchh hai ki  
coolie match jecta dekha hai  $\Rightarrow$  80% of the  
match.

Now lets plot a power law Distribution



80-20 rule



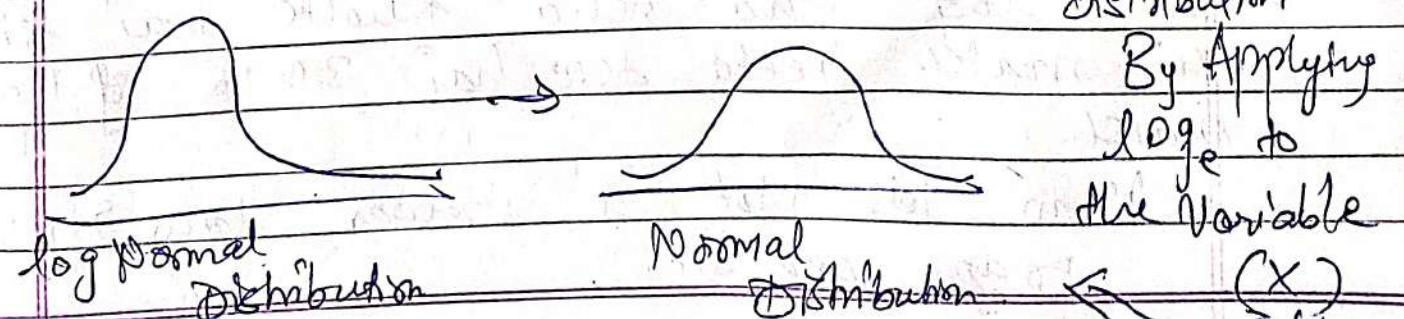
2<sup>nd</sup> example :- 80% of wealth is distributed with 20% of the total population

3<sup>rd</sup> Eg :- 20% of employees get higher salary with office than 80% of some Organisation employees -

4<sup>th</sup> Eg :- 80% of the expensive Cars are bought by 20% of people in the world.

5<sup>th</sup> Eg :- oil rich nation (Dubai, Saudi Arabia etc.)  
↳ 80% of the total oil is left with 20% of the nation.

Now, there is a question that previously we used to log Normal distribution  $\rightarrow$  Normal Distribution



Date \_\_\_\_\_

So, Can we convert this Power Law Distribution into Normal Distribution?

→ Yes!

We convert this power law Distribution into Normal Distribution Using Box Cox Transform.

will see in Coming Videos. ←

Now, whoever follows this type of Distribution we call it as Pareto Distribution.

And this Pareto Distribution Using the concept of power law Distribution distributes the Data.

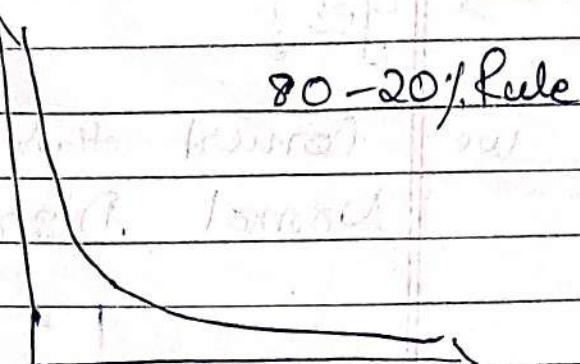
will see Pareto Distribution Math on later Videos.

But let understand that Pareto Distribution is not a Normal or Gaussian Distribution, It is a non-gaussian Distribution which follows Power law Distribution.

Whenever when we discuss about Pareto Distribution, we should always power law Distribution

To know about Power law Distribution

↳ Previous page.



Any Distribution which forms power law Distrib<sup>n</sup> then we specifically called it as pareto

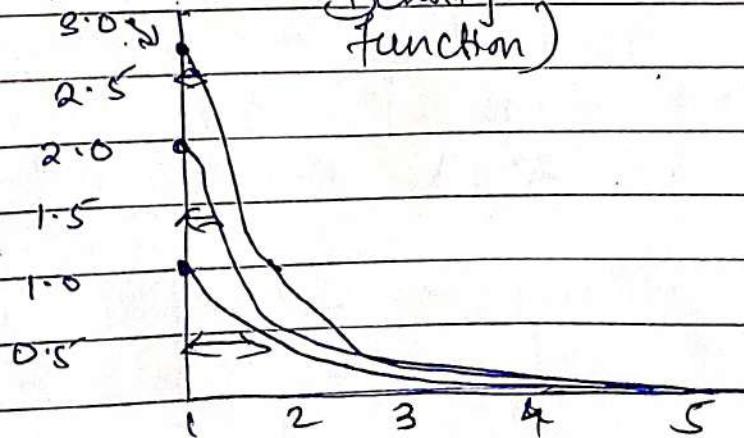
Interview Question

↳ What is Pareto Distribution?

Ans ⇒ It is an example of Power law Distribution and this Distribution are Non-Gaussian Distribution.

Let see some Examples :- Pareto Type 1

$\Pr(X=x)$  (Probability Density Function)



Alpha

$\alpha = \infty$

$\alpha = 3$

$\alpha = 2$

$\alpha = 1$

Alpha  
Decrease  
Height  
Decreases

If Alpha increases, the height also increases.  
and if you notice the gap between the blue and red gets reduced after  $\alpha$  increases.

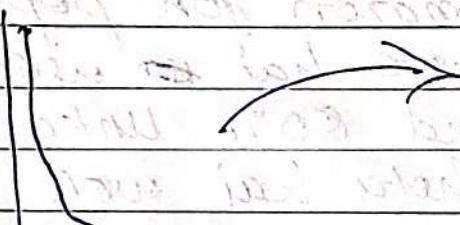
It follows Power law Distribution rule.

Specifically, if I talk about the height and the gap, then it is determined using  $\alpha$  (Alpha).

And why is it necessary to know, because if you train any model, like if I'm specifically working on Machine learning algorithms,

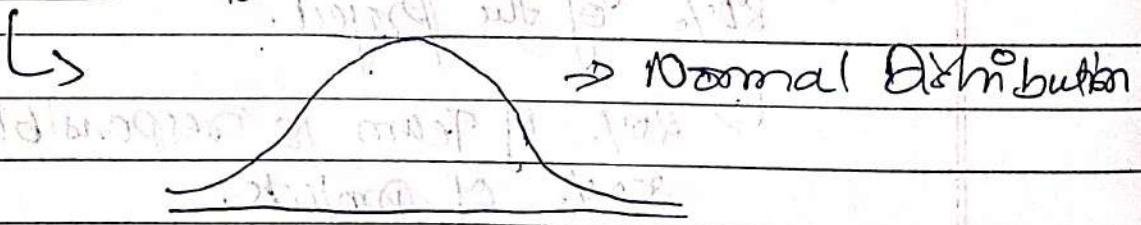
then some of the ML Algorithms will say that my data should be in normal distribution,

~~lets~~ Lets assume if I have a dataset which follows this distribution which specifically follows Power law Distribution)



Toh tya main  
ise ek normal  
Distribution main  
Change kar  
paunga

Can I change Power law Distribution into  
normal Distribution?



Note J

The importance of Normal Distribution is that any of your Dataset forms Normal Distribution then it gets easy to work with Machine Algorithms like Linear Regression, Logistic Regression we can work on such models easily and can create efficient models.

So) If my Dataset follows the 80-20% Rule of Pareto Distribution and I want to convert that into Normal distribution then the answer for this is Yes!!  
By Using Box Cox Transformation

↳ we'll see this using Python Programming.

### Pg: of Pareto Distribution

↳ 20% of the product in Amazon is responsible for 80% of Sales.

↳ (Amazon jok popular products hote hai ~~mech~~ with 20% hote hain And 80% unka sales revenue hote hain woh yani ee Aka hain)

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

③ 20% of Team is responsible for completing 80% of the Project.

↳ 20% of Team is responsible for delivering 80% of projects.

Hypothesis Testing

Date

It's very important topics for the people those who are approaching for the interview of Data Analyst, Data Science.

There is a huge chance of asking Question on 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.

↳ [wikipedia](#)

Another Defn:- Hypothesis Testing is an act 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.

↳ [Favstopeedia](#)

And in Both these Definition, we are talking about Inferential Statistics

In, Inferential Statistics, will probably take Sample Data. And from this Sample Data whatever insights or information we are getting, will come out with a conclusion about the Population Data.

To get the Conclusion from the Sample Data we use Hypothesis Testing!

And for this Hypothesis Testing, there are different techniques:-

Z - Test

T - Test

Anova Test

Chi Square Test.

Let's start with the steps of hypothesis testing.

Let's do an experiment on the coin, if it's fair or not.

→ So I took a coin, where head & tail is outcome, so I want to find out if this coin is fair or not.

Result  
How fair?

① So if we toss the coin and 50% of time I get head and another 50% of time I get tail then we can say that the coin is fair.

② Let say 60% of time - head

2 40% of time - tail

Then also this coin is fair. we can call 60 - 40 ratio. it as fair

~~60%~~ → ~~40%~~

③ 30% - head

70% - tail

then this coin is unfair

And ~~all~~ All this outcome we can find out through hypothesis theorem.

## Hypothesis Testing Steps :-

Page No.

Date

- Step ①  $\rightarrow$  NULL hypothesis ( $H_0$ ) = Coin is fair.
- ②  $\rightarrow$  Alternate Hypothesis ( $H_1$ ) = Coin is not fair
- ③ Experiment to prove ~~that~~ whether the NULL hypothesis should be Accepted or rejected.
- Experiment & Tossing coin for 100 times. sample of data
- ① 1st 100 times (Coin is fair or not)
- 50 - head (→ population)
- 50 - tail
- Coin is fair
- ② 2nd 100 times
- 60 - head
- 40 - tail
- Now, if the coin is fair we can consider the 60-40 ratio.
- ③ 3rd 100 times.
- 70 - head
- 30 - tail
- Now, here I get conscious that I got 70 times head and remaining 30 times tail. So By this we'll probably the coin is not fair.

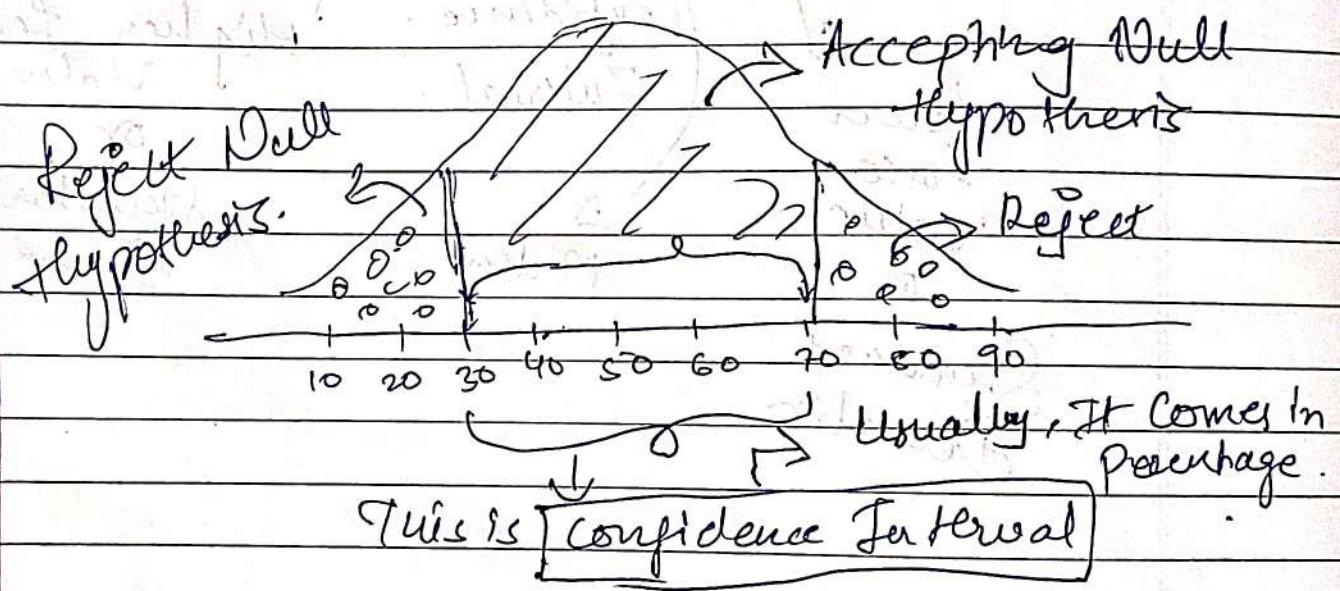
→ 2 Fair

Whenever we toss the coin, there are possibilities of going below 50 or Above 50 ~~as either~~

Let say the Average is 50 %.

Let Assume, if we get 90 times head, so ~~we~~ we obviously confirm that the coin is not fair.

And for that we deduce Confidence Interval



And how ~~to~~ Confidence Interval is defined by?

→ Confidence Interval is defined by Domain Expert.

Domain expert Pg: Let assume there is medical use case let say I want to take covid vaccination and I want to check whether this vaccination is right or not, so will see its impact and effect ~~as~~ and according to that ~~will~~ be the domain expert will decide its confidence Interval.

Confidence Interval is also called <sup>as Decision</sup> ~~as~~

Note:-

Page No. \_\_\_\_\_  
Date \_\_\_\_\_

Confidence Interval Usually Comes In  
Percentage.  $\rightarrow$

$$C.I = 95\%$$

lets take

C.I

2 Tail Test

Significance  
Value



95%

Significance  
Value

2.5%

Confidence -  
Interval.

Higher fence  
Value

or

Standard  
Deviation

Median

Lower  
fence  
Value

or

Standard  
Deviation