

## About Myself

Hi, I'm **Yash Pandey**, a passionate tech enthusiast From **Maharaja Agrasen Institute Of Technology (MAIT)** with a strong interest in **Artificial Intelligence (AI), Machine Learning (ML), and Data Science**. I have a deep curiosity for understanding how data-driven technologies shape the world, and I constantly strive to enhance my technical skills to stay ahead in this ever-evolving field.

My journey in the tech world has been driven by a problem-solving mindset and a desire to build impactful solutions. I enjoy working with complex datasets, developing predictive models, and optimizing machine learning algorithms. Apart from AI/ML, I have experience in **software development, competitive programming, and full-stack web development**. I believe in continuous learning and love sharing my knowledge with the tech community.

## My Skills

- **Machine Learning & AI:** Supervised & Unsupervised Learning, Deep Learning, Model Optimization
- **Web Development:** React.js, Node.js, Express.js, Database Management
- **Problem-Solving & DSA:** DSA In Java With 600+ LeetCode Problems .

## About These Notes

This document contains **comprehensive statistics notes** curated for the AIML learning path. It covers key statistical concepts essential for data analysis, machine learning, and AI applications. These notes are structured to be beginner-friendly yet detailed enough for deeper insights. They can serve as a reference for anyone looking to strengthen their foundation for AI/ML.

## Connect With Me

-  **GitHub:** [github.com/YashPandey1405](https://github.com/YashPandey1405)
-  **LinkedIn:** [linkedin.com/in/yashpandey29](https://linkedin.com/in/yashpandey29)
-  **LeetCode:** [leetcode.com/u/pandeyyash041](https://leetcode.com/u/pandeyyash041)

For **Statistics Code Implementation using Python**, check out my repository:

-  [GitHub - Statistics Using Python](https://GitHub - Statistics Using Python) 

- Statistics for AIML**
- Statistics is an Mathematical Science Including methods of Collecting, Organizing, Analysis of data in such a way That Meaningful Conclusion can be drawn from The Available data.
  - Data Is Facts / Piece of data That can be stored and measured.

- # **Uses of Statistics**
- Weather Forecasting : Prediction
  - Sports Analysis. SA : Interpret
  - Election Campaign.
  - E - Commerce platform.
  - Medical Industry.

→ **Types of Statistics**

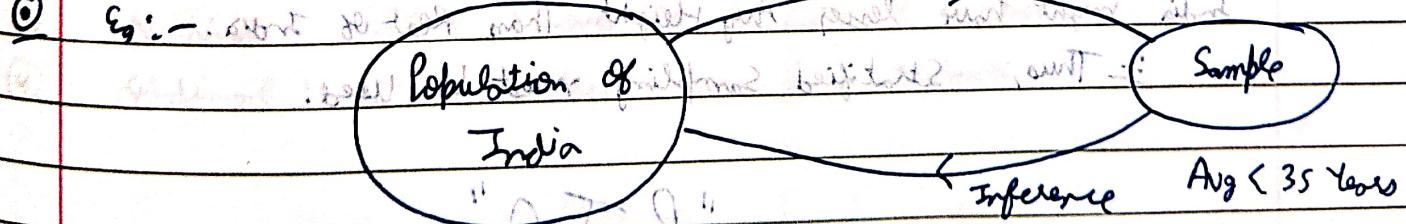
(1) Descriptive Statistics [Descriptive]

→ It Consist of Organizing and Summarize The data/Population.

Ex:- Virat Kohli Average In ODI's.

(2) Inferential Statistics [Inference]

→ It Consist of Using The data that has been measured To form The Conclusions about The Population.





#

## Types of Descriptive Statistics

- Measure of Central Tendency (Mean, Median, Mode)
- Measure of Symmetry (Skewness, Kurtosis)
- Measure of Dispersion

Note

→ Understand Descriptive VS Inferential by below :-

①

Descriptive : What Is The Avg. Height of The Class? etc.

②

Inferential : Are The Avg. Height of Students In Same what we expect from The Entire School.

→

Types of Sampling

1.

### Simple Random Sampling

→

Every Member of The Population ( $N$ ) has equal chance of Being Selected of your Sample.

②

Eg:- For Exit Poll, In an 12-Member Group, Each Person has Selection Probability of  $1/12$ .

2.

### Stratified Sampling [ Strata = layers / Groups ]

→

In This Sampling, different distinct categories are There.

→

An Simple Sample would be chosen from Each Strata (or layer).

③

Eg:- In Survey of Avg. Height In India, The people of North-East India might have lesser Avg Height Than Rest of India.

(Ans)

Thus, Stratified Sampling must be used.

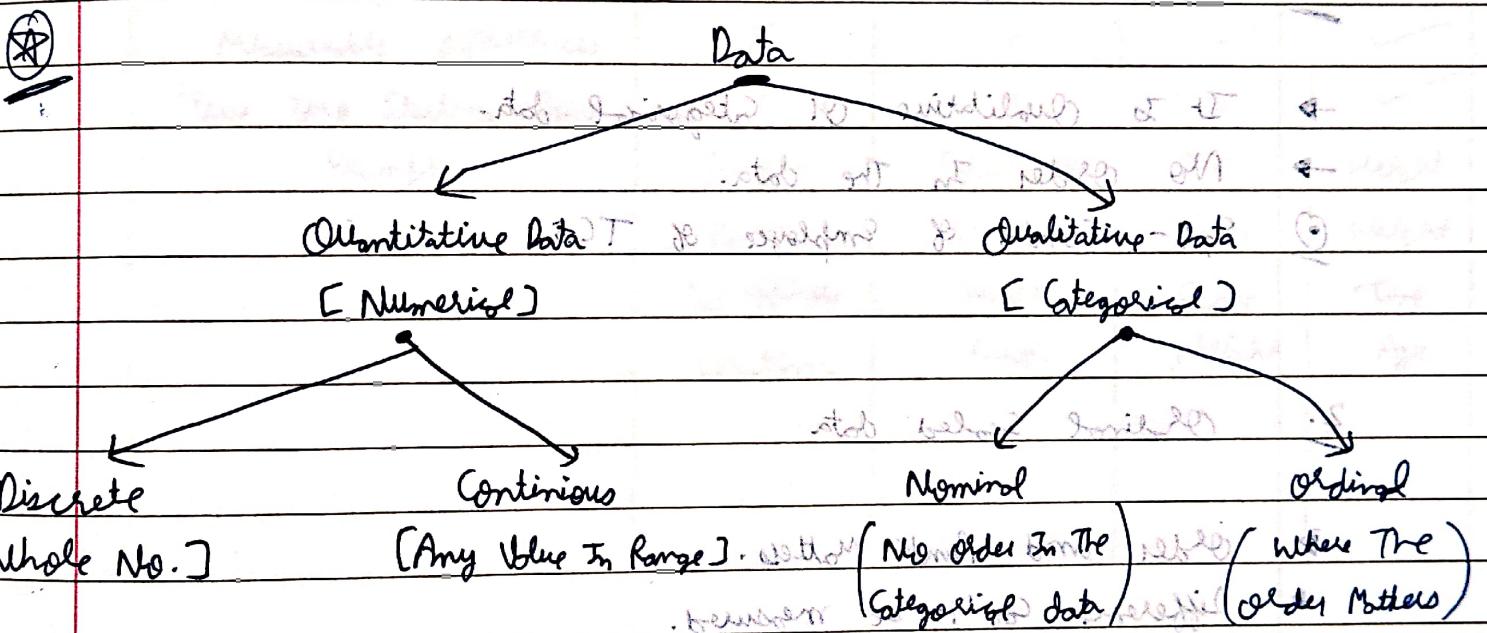
### ③ Cluster Sampling

- The division of The Sample Into Groups and Clusters where Some of The Clusters are been randomly selected.
- Then, all Individual in The chosen Cluster are Selected In The Sample.

### ④ Systematic Sampling

- Every  $n^{\text{th}}$  Element will be selected.

Ex:- Odd Roll Number.  
People Born In Odd Years.



① Discrete : N.o. of Childrens In a House.

② Continuous : Weight Of Students Of Class III.

③ Nominal : Blood Group Of Students Of School.

④ Ordinal : Marks Obtained By The Students.

## ⇒ Scales of Measurements

- ① Interval Scale (data doesn't have a starting point) } Quantitative Measures (Measurements are Numeric)
- ② Ratio Scale (data has a starting point) } Quantitative Measures (Measurements are Numeric)
- ③ Nominal Scale data } Qualitative data / Measures
- ④ Ordinal Scale data (It is classified in Non-Numeric) } Qualitative data / Measures

### 1. Nominal Scale data

- It is Qualitative or Categorical data.
- No order in the data.
- ① Eg:- Ranks of Employees of TCS [Qualitative]

### 2. Ordinal Scaled data

- Order (rank) Matters. (rank of 1st part)
- Difference can be measured.
- ① Eg:- Educational Qualification of Interns In HCL.

### 3. Interval Scaled Data

- The Rank and The order has no meaning.
- The difference can be measured by excluding Ratio.
- It doesn't have a starting value.
- ① Eg:- Height / Weight.

#### 4. Ratio Scaled data.

- The Order and Rank has an meaning.
- Differences and The Ratios are been measurable.
- It does have an '0' starting value.

Q) Eg :- Person A = 48 kg      }      A : B = 2 : 1  
                   Person 2 = 24 kg      }



Data labelled	Nominal	Ordinal	Interval	Ratio
Meaningful order	X	✓	✓	✓
Measurable differences	X	X	✓	✓
True zero Starting point	X	X	X	✓
Example	Gender Religion Post offices Location	Satisfaction Rating Mode Rank	IQ Temp Score Weight	Height Weight Time Age

## ⇒ Descriptive Statistics

- Summarization of data without Adding / Subtracting anything at any specific instance of time.
- 3 types are
  - Measures of Central Tendency.
  - Measures of Dispersion.
  - Measures of Symmetry.

### ⇒ Measures of Central Tendency

- Central Tendency represents the Centre Point of a dataset.
- It's 3 types are Mean, Median and Mode.

① Mean (Arithmetic mid-value of the data)

Mean =  $\frac{\sum x_i}{N}$

→ For Population,  $M_p = \frac{1}{N} \sum_{i=1}^N x_i$

→ For Sample,  $M_s = \frac{1}{n} \sum_{i=1}^n x_i$

Ex:- [1, 2, 3, 4, 5], mean =  $\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$

② Median (Physical Mid-Value of data)

→ If, Odd i.e. [1, 2, 3, 4, 5], Then, median = 3 (Mid Value)

→ If Even i.e. [1, 2, 3, 4, 5, 6], Then, median =  $\frac{3+4}{2} = 3.5$

★ Let, List = [1, 2, 3, 4, 5, 10000] → Here, 10000 is an outlier as it is very different from the rest of the

→ Outliers are the numbers which are much higher or lower than the other members. They are extreme values.

Note → Mean is affected by outliers, but not with median.

Case - I

$$[1, 2, 3, 4, 5]$$

mean = 3

median = 3

Case - II

$$[1, 2, 3, 4, 99]$$

mean = 21.8

median = 3

∴ Thus, Outliers affect the mean but not median.

- JAI LOGIC

③ Mode (The highest frequency element)

→ Let,  $[1, 2, 2, 3, 3, 3, 4] \rightarrow \text{mode} = 3$



Data

Continuous Data

Categorise data

Outlier is present

Outlier is not present

Replace / Input The Missing Value with The mode In The Column....

Median

Mean

(Replace The Missing Value with The Median of The Column)

(Replace The Missing Value with The Mean of The Column)

most used & easiest form of central tendency is the mean

## # Implementation of Central Tendency Using Python

→ import "numpy" as np # In The Program

① Mean → np.mean(list)

② Median → np.median(list)

③ Mode → Mode Is Not Present In The Numpy, for Mode,  
import statistics # Inbuilt Library In Python

— JAI LOHIE

(Thank You & Visit Again) BY



```
1 # Measures Of Central Tendency....  
2 import numpy as np # type: ignore  
3 import statistics  
4  
5 age=[15,75,29,86,45,58,37,41]  
6  
7 # Mean Calculation Using Numpy.....  
8 print(f"Mean Of {age} Is : {np.mean(age)}")  
9  
10 # Median Calculation Using Numpy.....  
11 print(f"Median Of {age} Is : {np.median(age)}")  
12  
13 # Mode Calculation Using Statistics.....  
14 print(f"Mode Of {age} Is : {statistics.mode(age)}")
```

⇒ Measures of Dispersion (and spread of a set of data in the given set) (i)

→ let,  $S_1 = \{1, 2, 3, 4, 5\}$ ,  $S_2 = \{3, 3, 3, 3, 3\}$

Mean = 3  $\times$  Median = 3  $=$  Mean = 3 (ii)  
Median = 3  $\times$  Median = 3

→ Spread of 2 data. → Range, IQR, Standard Deviation, Variance, etc.

minimum & maximum → Percentage and Percentile.

and 2nd quartile → Quartiles → to find Q1, Median, Q3

25.0 × 11 = 275.0 → (Mean) × 2F = 25.0 → (Variance) × 2F = 25.0

25.0 × 2.8 = 70.0 → Standard deviation = 2.8

(1) Range is difference b/w Max. Value & Min. Value

$$\{1, 2, 3, 4, 5\} \rightarrow \text{Range} = 4.8$$

$$\text{admit at P3} \{1, 2, 3, 4, 1000\} \rightarrow \text{Range} = 999$$

\* Outliers affect The Range.

(2) Percentage or Percentile

Outliers affect the range so it is given that without the outliers

→ let,  $S = \{1, 2, 3, 4, 5\}$ , Then,

$$S_1 = \% \text{ Numbers as odd} = \frac{3}{5} = 60\% \text{ numbers are odd}$$

and there are 5 even numbers

→ Percentile is an ~~value~~ value below which an certain percentage of the observations lies.

$$\text{Ex:- } S = \{1, 2, 3, 4, 4, 6, 7, 7, 8, 10\}$$

$$\rightarrow \text{Percentile Rank of } 6 = \frac{\# \text{ Values below } 6}{\# \text{ Values}} = \frac{5}{10} = 50^{\text{th}} \text{ Percentile}$$



$$\textcircled{1} \quad \text{Percentile of a Value} = \frac{\# \text{Values below The Value}}{\# \text{Values}} \times 100$$

$$\textcircled{2} \quad \text{Value} = \frac{\text{Percentile}}{100} \times (n+1) \quad \begin{matrix} \downarrow \\ \text{# Values} \end{matrix}$$

$$\textcircled{3} \quad \text{Ex:- } S = \{1, 2, 2, 2, 3, 4, 4, 4, 5, 6\}$$

. ~~arranged from smallest to largest~~

↑ 8<sup>th</sup> Number

Then, Value at 75<sup>th</sup> Percentile will be

$$\text{Value} = \frac{\text{Percentile} \times (n+1)}{100} = \frac{75 \times (10+1)}{100} = 8.25^{\text{th}} \text{ Number}$$

$\therefore 75^{\text{th}}$  Percentile = 8.25<sup>th</sup> Member of The Sample

~~which = 8<sup>th</sup> & 9<sup>th</sup> Number of The Sample~~ = 4<sup>th</sup> member

8.25<sup>th</sup> ~~member~~  $\rightarrow$  Take 8<sup>th</sup> Number.

$\therefore$  Thus, 8<sup>th</sup> ~~member~~  $\rightarrow$  Take Avg. of 8<sup>th</sup> & 9<sup>th</sup> Number

8.75<sup>th</sup>  $\rightarrow$  Take 9<sup>th</sup> Number.

~~so that all steps will be~~

### 3. Quartiles

~~arranged in ascending~~  $\textcircled{2}$

$\rightarrow$  Quartiles are Values That divide a list of Numbers Into Quarters.

$\textcircled{1}$  Put The Numbers In Order.

$\textcircled{2}$  Cut The Numbers Into 4 Equal Parts.

$\textcircled{3}$  The Quartiles Will be at The Cut.



Ex:-  $\{6, 8, 5, 5, 7, 3, 9\}$  Median value  $\Rightarrow 6$

Order =  $\{3, 5, 5, 6, 7, 8, 9\}$

Divide into  $\frac{1}{2} \leftarrow \frac{1}{2} \leftarrow (\text{odd}) \cdot 0$

Then,  $3 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9$

Divide  $\frac{1}{2.5} \leftarrow \frac{1}{2.5} \leftarrow 0$

Cut      Cut      Cut

Median - 1  $\leftarrow$  Median - 2  $\leftarrow$  Median - 3  $\rightarrow 0$

$(\theta_1) \quad (\theta_2) \quad (\theta_3)$

Divide  $\frac{1}{2.5} \leftarrow \frac{1}{2.5} \leftarrow 0$

# Imp. Trick for Jewellies  $\leftarrow 0 \leftarrow (\text{odd}) \cdot 0$

① When  $n = \text{odd}$   $\leftarrow 0 \leftarrow (2.5)$  ② When  $n = \text{even} - \text{odd}$

$$\rightarrow \theta_1 = \left( \frac{n+1}{4} \right)^{\text{th}} \quad \rightarrow \theta_1 = \frac{n}{4}^{\text{th}}$$

$$\rightarrow \theta_3 = \frac{3(n+1)}{4}^{\text{th}} \quad \rightarrow \theta_3 = \frac{3n}{4}^{\text{th}}$$

$$\rightarrow \theta_2 = \left( \frac{n+1}{2} \right)^{\text{th}} \quad \rightarrow \theta_2 = \text{Avg} \left( \frac{n}{2}^{\text{th}}, \left( \frac{n+1}{2} \right)^{\text{th}} \right)$$

$\rightarrow$  Let, Values  $\{1, 2, 3, 4, 5\}$ ,  $n = 6$

$$\text{Given } \theta_1 = \frac{n}{4}^{\text{th}} \quad n=6 \quad \left[ \frac{6}{4}^{\text{th}} = 1.5^{\text{th}} \right] = \text{Avg} \left( 1^{\text{th}}, 2^{\text{th}} \right) = \frac{1+2}{2} = 1.5$$

$$\therefore \theta_3 = \frac{3n}{4}^{\text{th}} = \frac{3 \times 6}{4}^{\text{th}} = 4.5^{\text{th}} = \text{Avg} (4^{\text{th}}, 5^{\text{th}}) = \frac{4+5}{2} = 4.5$$

$$\theta_2 = \text{Avg} \left( \frac{6}{2}^{\text{th}}, \frac{6+1}{2}^{\text{th}} \right) = \text{Avg} (3^{\text{th}}, 4^{\text{th}}) = \frac{3+4}{2} = 3.5$$

$\therefore$  Thus,  $\theta_1 = 1.5, \theta_2 = 3, \theta_3 = 4.5$ .

⇒ Five Point Summary {P, Q<sub>1</sub>, M, Q<sub>3</sub>, P, Q<sub>4</sub>} →

$$\{P, Q_1, M, Q_3, P, Q_4\} = \text{new set}$$

→ Q<sub>0</sub> (min) → 0% → 0<sup>th</sup> Percentile

$$P \quad | \quad F \quad | \quad M \quad | \quad Q_3 \quad | \quad E \quad \text{new set}$$

$$Q_1 \rightarrow 25\% \rightarrow 25^{\text{th}} \text{ Percentile}$$

$$Q_2 \rightarrow 50\% \rightarrow 50^{\text{th}} \text{ Percentile}$$

$$Q_3 \rightarrow 75\% \rightarrow 75^{\text{th}} \text{ Percentile}$$

$$Q_4 (\text{max}) \rightarrow 100\% \rightarrow 100^{\text{th}} \text{ Percentile}$$

⇒ Inter-Quartile Range (IQR)

$$P = 0 \quad | \quad Q_1 = 1 \quad | \quad M = 2 \quad | \quad Q_3 = 3 \quad | \quad P = 4$$

$$Q_0 = 0 \quad | \quad Q_1 = 1 \quad | \quad Q_2 = 2 \quad | \quad Q_3 = 3 \quad | \quad Q_4 = 4$$

$$\frac{M - Q_1}{2} = \frac{2 - 1}{2} = \frac{1}{2} \quad \frac{Q_3 - M}{2} = \frac{3 - 2}{2} = \frac{1}{2}$$

$$\text{IQR} = Q_3 - Q_1 = 3 - 1 = 2$$

∴ Outliers are the extreme values, Then,

$$\text{Lower fence} = Q_1 - (1.5 \times \text{IQR})$$

$$\text{Upper fence} = Q_3 + (1.5 \times \text{IQR})$$

∴ Thus, Values  $E = [\text{Lower fence}, \text{Upper fence}]^T$  are valid values and the values outside this range are outliers.

$$P_{1.5} = \frac{(M - Q_1)}{2} = \frac{2 - 1}{2} = \frac{1}{2}$$

## " ≠ Box - Whisker Plot "

Q → Let, Values = {2, 3, 3, 3, 3, 4, 4, 5, 5, 5, 6, 6, 6, 7, 8, 99}  $(n=16)$

$$[2, 3, 3, 3, 3, 4, 4, 5, 5, 5, 6, 6, 6, 7, 8, 99] = \text{values}$$

$$\rightarrow Q_1 = \text{min value} = 0^{\text{th}} \text{ percentile} = 2.5+1 = 3$$

$$Q_1 = 25^{\text{th}} \text{ percentile} = \frac{25}{100} \times 17^{\text{th}} = 4.25 = 4^{\text{th}} = 3$$

$$Q_2 = 50^{\text{th}} \text{ percentile} = \frac{50}{100} \times 17^{\text{th}} = 8.5^{\text{th}} = 5+5 = 5$$

$$Q_3 = 75^{\text{th}} \text{ percentile} = \frac{75}{100} \times 17^{\text{th}} = 12.75^{\text{th}} = 13^{\text{th}} = 6$$

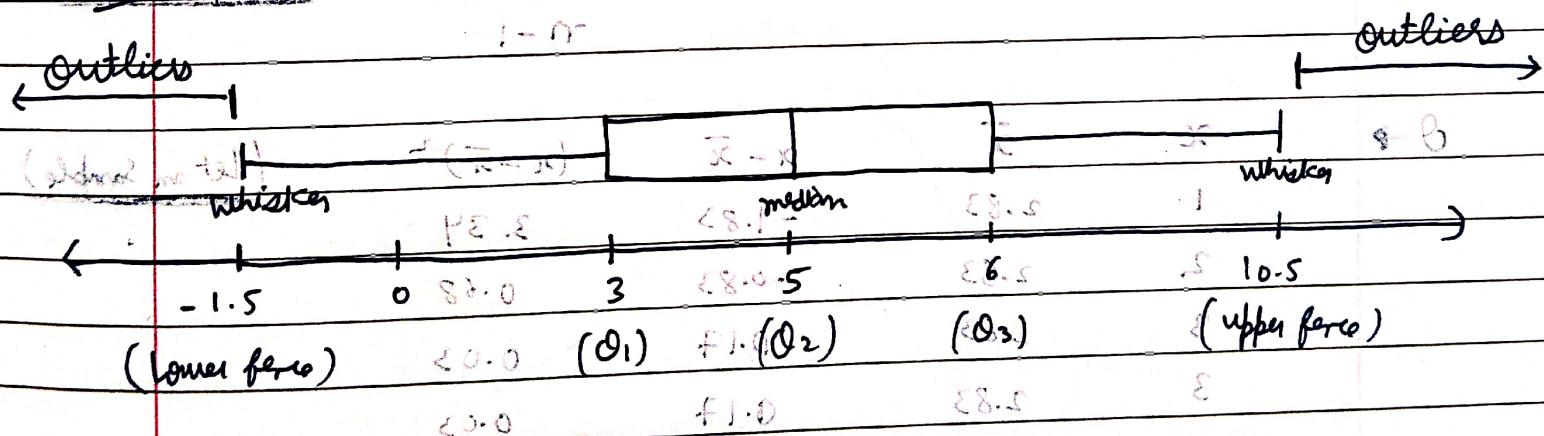
$$Q_4 = \text{max value} = 100^{\text{th}} = 99$$

$$\therefore \text{IQR} = Q_3 - Q_1 = 6 - 3 = 3 : \text{answ}$$

$$\text{(Lower fence)} = Q_1 - (1.5 \times \text{IQR}) = 3 - (1.5 \times 3) = 3 - 4.5 = -1.5 \quad (i)$$

$$\text{Upper fence} = Q_3 + (1.5 \times \text{IQR}) = 6 + (1.5 \times 3) = 6 + 4.5 = 10.5$$

The Box - Whisker Plot will be :- answ graph (ii)



$\therefore$  Thus,  $99$  value is outlier for given range of values.

" Well practice = good habit "

★ → Mean = Deviation  $\Rightarrow$  Average of  $(x - \bar{x})$  + mean  $= \bar{x}$

Q:- Let, values = [1, 2, 3, 4, 5]

$$\text{Mean} = \frac{1+2+3+4+5}{5} = 3 \text{ (mean) } = 10$$

$$\Sigma = 15 = 5 \times 3 = \text{Number of terms} = 5$$

$$\text{Mean Dev.} = |1-3| + |2-3| + |3-3| + |4-3| + |5-3|$$

$$= 2+2 = 4 = 1 \times 5 = \text{Number of terms} = 5$$

$$\text{Mean Deviation} = \frac{2+1+0+1+2}{5} = \frac{6}{5} = 1.2$$

$$\delta = \sigma = \sqrt{2^2 + 1^2 + 0^2 + 1^2 + 2^2} = \sqrt{10} = \text{Standard Deviation} = 3.16$$

∴ On an Average, Each of The data Point are 1.2 Unit away from mean value  
 $\therefore \text{PP} = 100 = \text{standard deviation} = 3.16$

4. Variance : The Average of Squared differences from the mean.

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

$$2.01 = (2^2 + 1^2 + 0^2 + 1^2 + 2^2) / 5 = (4+1+0+1+4) / 5 = 10 / 5 = 2.0$$

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

$\theta \rightarrow$	$x$	$\bar{x}$	$x - \bar{x}$	$(x - \bar{x})^2$	$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ (Let on Sample)
1	2.83	2.83	-0.83	0.69	0.69
2	2.83	2.83	-0.83	0.69	0.69
3	2.83	2.83	-0.17	0.03	0.03
3	2.83	2.83	0.17	0.03	0.03
4	2.83	2.83	2.17	4.70	4.70
4	2.83	2.83	1.17	1.37	1.37
				6.82	6.82

$$\therefore s^2 = \frac{\sum (x - \bar{x})^2}{n-1} = \frac{6.82}{5} = 1.37$$



## # How To Calculate Variance

- ① Calculate The Mean.  $\bar{x} = \frac{\sum x}{n}$ , where  $n$
- ② For Each No. In data, subtract The mean from Number.
- ③ Square The differences.  $(\bar{x}-x)$
- ④ Calculate The Average of Square of differences.

Note → When  $\sigma^2 \uparrow$ , Spread ↑ & when  $\sigma^2 \downarrow$ , Spread ↓

$$\text{Variance} = \frac{\sum (\bar{x}-x)^2}{n} > \frac{\sum (\bar{x}-x)^2}{1-n}$$

$\bar{x}$  >  $1-n$

5. Standard Deviation : The Measure of How spread Numbers are :-

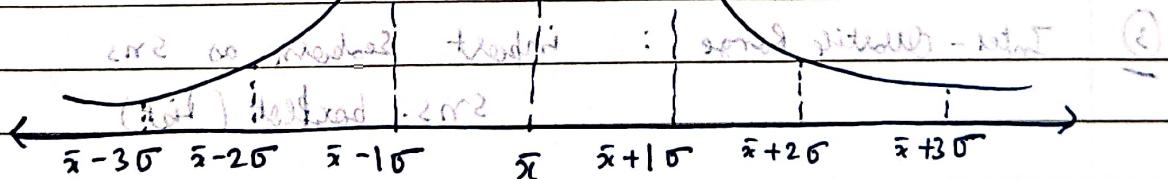
$$\therefore S_{\text{pop}} = \sqrt{\sigma^2_{\text{population}}} = \sqrt{\text{Variance}} \approx \sqrt{1.137} = \underline{\underline{1.17}}$$

$$\sigma_p = \sqrt{\sigma^2_{\text{population}}} \quad (\text{for a population}) \quad \text{and}$$

$$\sigma_s = \sqrt{\sigma^2_s} \quad (\text{Sample})$$

(Total) range - (Total) spread : spread (1)

(Ex.) [25, 102, 25, 6] (Total) spread / 4 : spread (2)



1 std. deviation : spread (3)

Q. Eg:- The Range for 1 standard deviation will be :- (4)

$$(2.83 - 1.17, 2.83 + 1.17) = [1.66, 4]$$

All Points In This Range are

within at 1 st. dev. from Mean ...

Note → In Population,  $\sigma_p^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$

In Sample,  $\sigma_s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)}$

$\therefore$  Now,  $\sum (x_i - \mu)^2 \gg \sum (x_i - \bar{x})^2$

$\therefore$  Thus, we use den. of Sample Smaller To make Ratios Equal. Thus, Instead of  $(n)$ , we use  $(n-1)$  In The denominator.

$$\frac{\sum (x_i - \bar{x})^2}{n-1} < \frac{\sum (x_i - \mu)^2}{N}$$

## # 5.1. Implementation of Method of Dispersion In Python :-

→ Use ( Numpy as np ) In the Program.

① Range :  $\max(\text{List}) - \min(\text{List})$

② Percentile :  $\text{np.percentile}(\text{List}, [0, 25, 50, 75, 100])$

③ Inter-Quartile Range : import Seaborn as sns

sns.boxplot(List)

④ Variance :  $\text{np.var}(\text{List})$

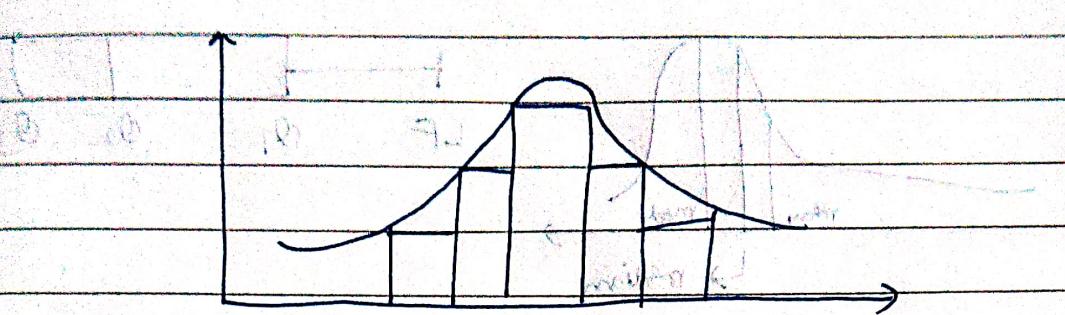
⑤ St. dev. :  $\text{np.std}(\text{List})$



```
1 import numpy as np # type: ignore
2 import seaborn as sns # type: ignore
3 import statistics
4
5 Numbers=[2,3,4,5,7,8,9]
6
7 # Range
8 print(f"The Range Will Be : {max(Numbers)-min(Numbers)}")
9
10 # Percentile
11 ans=np.percentile(Numbers,[0,25,50,75,100])
12 print(f"The Percentile Will Be : {ans}")
13
14 # Inter-Quartile Range(ITR) -> Q3 - Q1...
15 # Lower-Fence -> Q1 - (1.5*IQR)
16 # Upper-Fence -> Q3 + (1.5*IQR)
17 print(sns.boxplot(Numbers))
18
19 # Variance
20 print(f"The Variance By Numpy Will Be : {np.var(Numbers)}")
21 print(f"The Sample-Variance By Statistics Will Be : {statistics.variance(Numbers)}")
22 print(f"The Public-Variance By Statistics Will Be : {statistics.pvariance(Number
23 s)}")
24 # Standard-Deviation
25 print(f"The Standard-Deviation Will Be : {np.std(Numbers)})
```

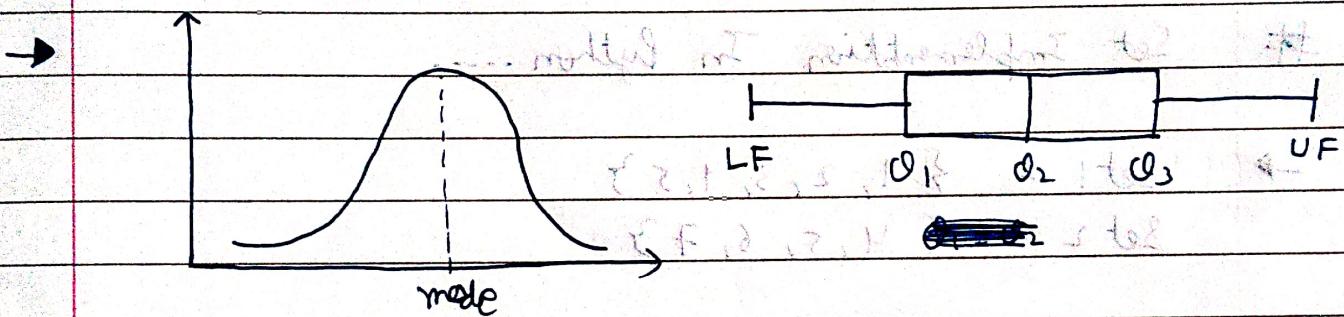
→ Measure of Symmetry

→ If Anything is Exactly Towards The left and Right.



→ Skewness : measure of data symmetry

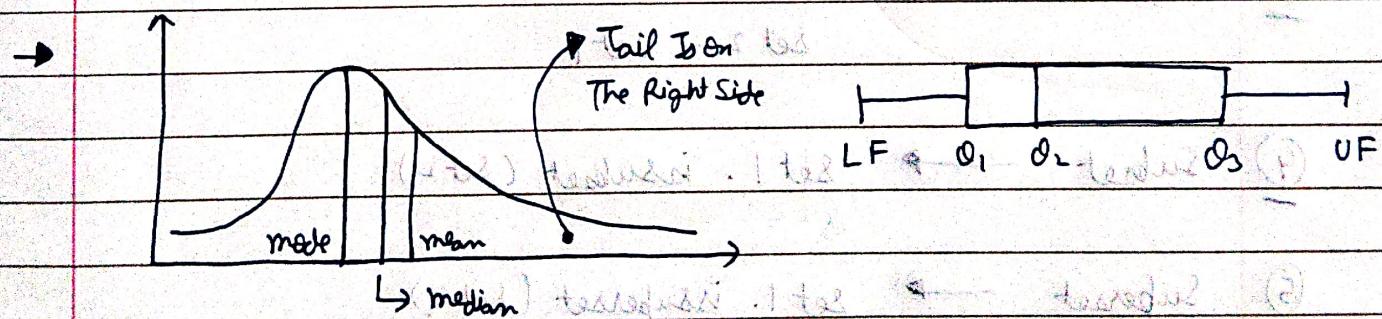
① No Skewness (Skewness = 0)



→ In No | Skewness,  $(\text{mean} = \text{median} = \text{mode})$

$$Q_2 - Q_1 \approx Q_3 - Q_2$$

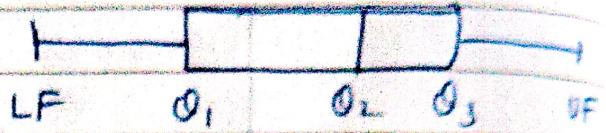
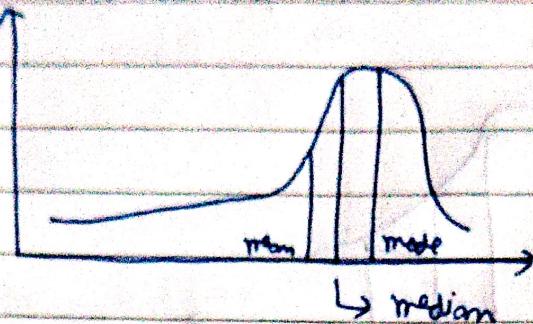
② +ve Skewed / Right Skewed (Skewness = +)



→ In (+ve Skewness),  $\text{mode} < \text{median} < \text{mean}$

$$Q_2 - Q_1 \ll Q_3 - Q_2$$

③

-Ve Skewed / Left Skewed (Skewness = -)

In - Ve Skewed, mean &lt; median &lt; mode

$$O_1 - O_2 \geq O_3 - O_2$$

#

Set Implementation In python ---



$$\text{Set 1} = \{1, 2, 3, 4, 5\}$$

$$\text{Set 2} = \{4, 5, 6, 7\}$$

①Union  $\rightarrow$  Set1.union(Set2), (set1 | set2)②Intersection  $\rightarrow$  Set1.intersection(Set2), (set1 & set2)③Difference  $\rightarrow$  Set1 - Set2

Set2 - Set1

④Subset  $\rightarrow$  Set1.issubset(Set2)⑤Superset  $\rightarrow$  Set1.issuperset(Set2)⑥Symm. Diff.  $\rightarrow$  Set1.symmetric\_difference(Set2), (set1 ^ set2)



```
1 set1 = {1,2,3,4,5}
2 set2 = {4,5,6,7,8}
3
4 # Union Of 2 Sets -> union = set1 | set2
5 print("Union of set1 and set2: ", set1.union(set2))
6
7 # Intersection of 2 sets -> -> intersection = set1 & set2
8 print("Intersection of set1 and set2: ", set1.intersection(set2))
9
10 # Difference of 2 sets -> diff = set1 - set2
11 print("Intersection of set1 and set2: ", set1.difference(set2))
12
13 # Subset And Superset In Python Sets...
14 print("Subset of set1 and set2: ", set1.issubset(set2))
15 print("Superset of set1 and set2: ", set1.issuperset(set2))
16
17 # Symmetric Difference Of 2 Sets -> SymmDiff = set1 ^ set2
18 print("Symmetric Difference of set1 and set2: ", set1.symmetric_difference(set2))
```

## ⇒ Covariance and Correlation

→ There are 2 Types of Relation as:-

- ① Direct  $\rightarrow X \propto Y$  (Ex:- House Area  $\propto$  Bigg Price)
- ② Indirect  $\rightarrow X \propto \frac{1}{Y}$  (Ex:- Age  $\propto \frac{1}{Health}$ )

### ① Covariance

$$\rightarrow COV(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

→ Variance Was spread of data, whereas, Covariance means you are been trying to understand The Relationship of one feature w.r.t others.

$$\theta \rightarrow \text{Let, } X = [2, 3, 6, 1] \rightarrow \bar{x} = 3$$

$$Y = [3, 5, 6, 8] \rightarrow \bar{y} = 5.5$$

X	Y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$
2	3	-1	-2.5	2.5
3	5	0	-0.5	0
6	6	3	0.5	1.5
1	8	-2	2.5	-5
				-1

$$\therefore COV(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1} = \frac{-1}{4-1} = \frac{-1}{3} = -0.34$$

∴ Thus, x and y are  $\Theta$ ve ~~not~~ Related To Each other.

∴ Thus,  $(X \propto \frac{1}{Y})$ . Made with Love by Yash Pandey

$\Theta \rightarrow$  Let,  $X = [2, 4, 6] \rightarrow \bar{x} = 4$   
 $Y = [3, 5, 7] \rightarrow \bar{y} = 5$

$X$	$Y$	$X - \bar{x}$	$Y - \bar{y}$	$(X - \bar{x})(Y - \bar{y})$
2	3	-2	-2	4
4	5	0	0	0
6	7	2	2	4
				8

$\therefore \text{Cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{(n-1)} = \frac{8}{3-1} = \frac{8}{2} = 4 \equiv \oplus$

$\therefore$  Thus,  $x$  and  $y$  are  $\oplus$ ve Related i.e.  $X \propto Y$ .

### ~~# Disadvantage of Covariance~~

- ① No Comparison of The Strength In The Covariance.
- ② No any Standardization Scale To Interpret The Strength.
- ③ Covariance has dimensions, which makes The Comparison more difficult.

$$(P-E)(X-\bar{x}) = P-E \quad X-\bar{x} \quad Y \quad X$$

$$7.5 \quad 7.5 \quad 1 \quad 8 \quad 5$$

$\Theta :-$  Let,  $\text{Cov}(\text{Height}, \text{Weight}) = 0.24 \text{ m.kg}$  (1)

$\text{Cov}(\text{Height}, \text{Salary}) = 54 \text{ m.Rs}$  (2)

$\text{Cov}(\text{Weight}, \text{Salary}) = 148 \text{ kg.Rs}$  (3)

$\rightarrow$  No Comparison of data Possible with ①, ② & ③.

Q2 Pearson Correlation Coeff.  $[-1, 1]$

$$\rightarrow \rho_{x,y} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = [-1 \text{ to } 1]$$

$$\rightarrow -1 = \text{No Correlation} \quad 0 = (\text{P}(0)(\bar{x}-\bar{y})^2) = (\bar{P}_{xy})^2 = 1$$

The more 0, The more

-ve Correlated feature

The more +, The more

+ve Correlated feature

Q3: -  $\rho_{x,y} = 0.4$  } Feature (A, B) are Highly Correlated  
 $\rho_{A,B} = 0.8$  } As Compared To Feature (x, y) ---.

Q4 Calculate  $\rho_{x,y}$ ,  $x = [2, 3, 6, 1] \rightarrow \bar{x} = 3$   
 $y = [3, 5, 6, 8] \rightarrow \bar{y} = 5.5$

$x$	$y$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$
2	3	-1	-2.5	-2.5
3	5	0	-0.5	0
6	6	3	0.5	+1.5
1	8	(1-2)(8-5.5) = -2.5	= -2.5	-5
				-1

$$\underline{Q5} \quad \sigma_x^2 = \frac{\sum (x - \bar{x})^2}{n-1} = \frac{(-1)^2 + 0^2 + 3^2 + (-2)^2}{4-1}$$

$$\sigma_x^2 = \frac{1+0+9+4}{3} = \frac{14}{3} = 4.67$$

$$\therefore \sigma_x = \sqrt{\sigma_x^2} = \sqrt{4.67} = 2.161$$



$$\textcircled{2} \quad \sigma^2_y = \frac{\sum (y - \bar{y})^2}{n-1} = \frac{(-2.5)^2 + (-0.5)^2 + (0.5)^2 + (2.5)^2}{4-1}$$

$$= \frac{6.25 + 0.25 + 0.25 + 6.25}{3} = \frac{13}{3} = 4.33$$

$$\therefore \sigma_y = \sqrt{\sigma^2_y} = \sqrt{4.33} = 2.08$$

$$\textcircled{3} \quad \text{Cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1} = \frac{-1}{4-1} = \frac{-1}{3} = -0.33$$

$$\therefore \rho_{x,y} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{-0.33}{2.16 \times 2.08} = -0.0734$$

~~both are not linear so (x,y) and (y,x)~~

$\therefore$  Thus,  $x$  is very less correlated with  $y$ .

Note → Pearson Correlation Coeff. always measures The Linear Relationship.

$$2.2 < 5 \quad \therefore (2.2, 2.8) \neq y$$

→ For Non-Linear Relationship, Use Spearman Rank Correlation

$$(P-P) \quad (x-x) \quad (y-y) \quad (x-y) \quad y \quad x$$

$$2.3 \quad 2.5 \quad 2.8 \quad 2.2 \quad 2 \quad 5$$

⇒ Spearman Rank Correlation ( $y$ )

$$1.1 \quad 2.0 \quad 3.0 \quad 2.0 \quad 2 \quad 5$$

→  $\gamma_{x,y} = \frac{\text{Cov}(R(x), R(y))}{\sigma_{R(x)} \cdot \sigma_{R(y)}}$

$$1.1 \quad 2.0 \quad 3.0 \quad 2.0 \quad 2 \quad 5$$

① where,  $R(x) = \text{Rank of } x$

$R(y) = \text{Rank of } y$

$$\rightarrow \text{Find } Y_{X,Y}, X = [5, 7, 8, 1, 2] \rightarrow \bar{x} = 4.6 \\ Y = [6, 4, 3, 1, 2] \rightarrow \bar{y} = 3$$

$$\rightarrow \text{For } R(X), \text{ Sort } X \rightarrow X = [1, 2, 5, 7, 8] \\ \text{⑤ ④ ③ ② ①}$$

(Numbers In Reverse Order)

$$\rightarrow \text{For } R(Y), \text{ Sort } Y \rightarrow Y = [1, 2, 3, 4, 6] \\ \text{⑤ ④ ③ ② ①}$$

$$\therefore P(X) - (R(X)) \cdot (R(Y)) = (R(Y)) - (R(X))$$

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

$$\therefore \text{Thus, H.O. } R(X) = [3, 2, 1, 5, 4] \rightarrow \overline{R(X)} = 3.$$

$$R(Y) = [1, 2, 3, 5, 4] \rightarrow \overline{R(Y)} = 3$$

$$\therefore R(X) \text{ and } R(Y) \text{ have } R(X) - \overline{R(X)} \text{ and } R(Y) - \overline{R(Y)} \text{ as Product}$$

3	0	-2
2	-1	-1
1	-2	0
5	2	4

$$\textcircled{1} \sigma^2_{R(X)} = \frac{\sum (R(X) - \overline{R(X)})^2}{n-1} = \frac{0^2 + (-1)^2 + (-2)^2 + 2^2 + 1^2}{5-1}$$

$$\sigma^2_{R(X)} = \frac{0+1+4+4+1}{4} = \frac{10}{4} = 2.5$$

∴ Thus,  $\sigma_{R(X)} = \underline{\text{Made with Love by Yash Pandey}}$



$$\textcircled{2} \quad \sigma^2_{R(Y)} = \frac{\sum (R(Y) - \bar{R}(Y))^2}{n-1} = \frac{(-2)^2 + (-1)^2 + 0^2 + 2^2 + 1^2}{5-1}$$

$$\sigma^2_{R(Y)} = \frac{4+1+0+4+1}{4} = \frac{10}{4} = 2.5$$

$$\therefore \sigma_{R(Y)} = \sqrt{\sigma^2_{R(Y)}} = \sqrt{2.5} = 1.581$$

$$\textcircled{3} \quad \text{Cov}(R(X), R(Y)) = \frac{\sum (R(X) - \bar{R}(X))(R(Y) - \bar{R}(Y))}{n-1} = \frac{4}{5-1} = \underline{\underline{1}}$$

$$\therefore \rho_{R(X), R(Y)} = \frac{\text{Cov}(R(X), R(Y))}{\sigma_{R(X)} \cdot \sigma_{R(Y)}} = \frac{1}{1.581 \times 1.581} = \underline{\underline{0.4}}$$

$\therefore$  Thus,  $\rho_{R(X), R(Y)} = 0.4$

$\therefore$  Thus, Spearman Rank Correlation  $(X, Y)$  Is 0.4.

Note → If we have to try Combinations which can result in better growth of the company, we can use correlation there.

③ Let, Product =  $X$ , Options =  $\{A, B, C\}$

→ Let, On computation from data,  $\rho_{X,A} = -0.8 \rightarrow$  Neg. Reln  
 $\rho_{X,B} = 0.01 \rightarrow$  No Reln  
 $\rho_{X,C} = 0.9 \rightarrow$  Pos. Reln

$\therefore$  Thus, we can use "X-C" Combo here.

```
● ● ●

1 import pandas as pd
2 import numpy as np
3
4 # Load your DataFrame (Example data, replace with actual dataset)
5 df = pd.DataFrame({
6     'A': [1, 2, 3, 4, 5],
7     'B': [2, 3, 4, 5, 6],
8     'C': ['Yes', 'No', 'Yes', 'No', 'Yes'] # Non-numeric column
9 })
10
11 # Convert categorical values ('Yes'/'No') to numeric if present
12 df.replace({'Yes': 1, 'No': 0}, inplace=True)
13
14 # Select only numeric columns
15 df_numeric = df.select_dtypes(include=[np.number])
16
17 # Calculate the correlation matrix
18 corr_matrix = df_numeric.corr()
19
20 # Print the correlation matrix
21 print("Correlation Matrix:")
22 print(corr_matrix)
23
```



## # Implementation of Cov (x,y) and $S_{x,y}$ In Python

- ① Load The data In The Variable Let "df" Is A Pandas DataFrame
- ② For Covariance :  $df.cov()$
- For Correlation :  $df.corr()$

- JAI LOGIC

⇒ Random Variable

→ A Set of Possible Values from an Random Experiment.

① Eg:- Let, Sample = { HHH, HHT, HTH, THH, HTT, THT, TTH, TTT }  
 $X = \# Heads$

$$\therefore \text{Then, } P(X=1) = 3/8$$

$$P(X=2) = 3/8$$

$$P(X=3) = 1/8$$

$$P(X=0) = 1/8$$

- Random Variable Is an Set of possible Values from an Random Exp.
- An Random Variable Value Is been Unknown.
- An function which assigns Values To Each of Exp. outcomes.

→ Probability dist<sup>n</sup>  $F^n$  (PDF)

- Irrespective of outcome "Nature", draw outcomes in an form of distribution. This is known as Prob. dist<sup>n</sup>  $F^n$ .
- The 2 types are
  - Discrete (Prob. Mass  $F^n$ )
  - Continuous (Prob. density  $F^n$ )

QUESTION 1 :-

### ① Probability Mass Function

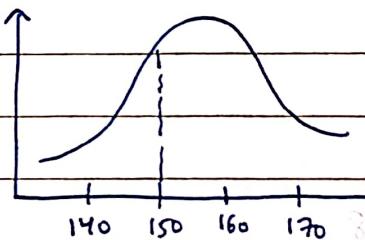
→ The Distribution of The discrete Random Variable.

② Eg:- Rolling a dice. No. most likely to get A

### ② Probability Density Function

→ The Distribution of The Continuous Data.

③ Eg:-



$$P(X = 155) = 0.1$$

$$P(X \leq 150) = ?$$

$$P(X \leq 150) = P(Z \leq -0.5)$$

$$P(Z \leq -0.5) = 0.3085$$

### ③ Cumulative Distribution Function (CDF)

→ CDF Is The Summation of all Probabilities Possible upto an Given Point.

④ Eg:- Let The Case for Rolling of dice:-

X	P(X)	CDF
1	$\frac{1}{6}$	$\frac{1}{6}$
2	$\frac{1}{6}$	$\frac{2}{6}$
3	$\frac{1}{6}$	$\frac{3}{6}$
4	$\frac{1}{6}$	$\frac{4}{6}$
5	$\frac{1}{6}$	$\frac{5}{6}$
6	$\frac{1}{6}$	1

## # Different Types of Probability Distribution $F^n$ (PDF) (A29, 249)

→ The 2 types of probability distribution  $F^n$  are :-

(i) ~~Cont.~~ (Continuous) → ~~Assume an interval will help in finding out~~

Normal Dist<sup>n</sup> / Gaussian Dist<sup>n</sup>.

$\Rightarrow Z = (X - \mu) / \sigma$  → Standard Normal Dist<sup>n</sup>:  $X = ?$

2. Prob. density F<sup>n</sup>

→ Log-Normal Dist<sup>n</sup>.

→ Chi-Square Dist<sup>n</sup>.

→ F-Distribution:  $F = ?$

(discrete)  $P = ?$

1. Prob. Mass F<sup>n</sup>

→ Bernoulli Dist<sup>n</sup>.

→ Binomial Dist<sup>n</sup>.

→ Poisson Dist<sup>n</sup>.

★ Uniform Dist<sup>n</sup> Is Both Discrete & Continuous PDF.

→ Discrete Uniform Distributions → DISCRETE

→ An Uniform Distribution refers To The Type of an Probability Dist<sup>n</sup> in which The outcomes are equally likely.

→ Uniform Dist<sup>n</sup> are Both Continuous and Discrete.

→ Notation of Uniform Dist<sup>n</sup> =  $U(a, b)$  �  $\{a, b\}$

→ Mean of Discrete Uniform Dist<sup>n</sup> =  $\frac{a+b}{2}$ ,  $a = \text{start}$ ,  $b = \text{end}$

Variance of Discrete Uniform Dist<sup>n</sup> =  $\frac{n^2 - 1}{12}$

H.S.D =  $H.o.X.o = (a-b)^2 / 12$



**DISCRETE**

→ Bernoulli Dist<sup>n</sup> (Binary  $\rightarrow$  2 outcomes)

→ An Discrete Prob. Distribution of a Random Variable which Takes only Two Possible Outcomes, Typically labelled as Success (Coded As 1), Failure (Coded as 0) with a fixed Prob. of Success & Failure.

Q :-  $X = \text{No. of Heads in one Toss}$   $\rightarrow$  Success (Head) = 0.5 =  $p$   
Failure (Tail) = 0.5 =  $1-p$

$$\therefore \text{bmf} = P(X=k) = \begin{cases} p & \text{if } k=1 \text{ (Success)} \\ 1-p & \text{if } k=0 \text{ (Failure)} \end{cases}$$

$$P(X=k) = p^k(1-p)^{1-k}$$

## # Conditions of Bernoulli Dist<sup>n</sup>

- No. of Trials must be finite.
- Each Trial must be Independent.
- There are only 2 possible outcomes as Success & Failure.
- Prob. of each outcome should be same in every trial.



In Bernoulli Dist<sup>n</sup>, Mean =  $p$  and Variance =  $p(1-p)$

Q → Bumrah Bowls 6 Balls at wicket with the Prob. of 0.6 at hitting the Stumps with each ball. What is the Prob. of Not Hitting in wicket?

→ Prob. of Hitting (Success) =  $p = 0.6$   
Prob. of Not-Hitting (Failure) =  $1-p = 0.4$

Q →  $\therefore$  Mean =  $p = 0.6$   
Variance =  $p(1-p) = 0.6 \times 0.4 = 0.24$

Q → St. Dev =  $\sqrt{\text{Var}} = \sqrt{0.24} = 0.4899$

~~DISCRETE~~

Binomial Dist<sup>n</sup> (Bernoulli with n-Trials)

~~CONTINUOUS~~

→ Binomial Dist<sup>n</sup>. If The n trials Bernoulli Trials are independent then

$$pmf = {}^n C_k \cdot p^k \cdot (1-p)^{n-k}$$

# Imp. Points about Binomial Dist<sup>n</sup>

→ The Binomial dist<sup>n</sup> Is PdF i.e.  $\sum p(x) = 1$ .

→ Moment Generating F<sup>n</sup> (MGF) =  $(pe^t + q)^n$

→ Characteristic F<sup>n</sup> (Chg F<sup>n</sup>) =  $(Pe^{it} + q)^n$

→ Prob. Generating F<sup>n</sup> (PGF) =  $(pz + q)^n$

→ In Binomial Dist<sup>n</sup>, Mean =  $np$  and Variance =  $npq$

$$(Mean)^2 = (np)^2$$

Q → With 3 Tosses, what is the prob. of getting exactly 2 Heads?

(Ans: 0.375)  $\Rightarrow$   $n=3, k=2, p(\text{Head}) = p = 0.5, q = 1-p = 0.5$

→ Here,  $n=3, k=2, p(\text{Head}) = p = 0.5, q = 1-p = 0.5$

$$\therefore ① P(X=2) = {}^n C_k \cdot p^k \cdot q^{n-k}$$

$$= 3 C_2 (0.5)^2 (0.5)^{3-2}$$

$$= \frac{3!}{2!1!} \cdot \left(\frac{1}{2}\right)^2 \cdot \left(\frac{1}{2}\right)^1$$

$$\therefore ② = \frac{3 \times 2 \times 1}{2 \times 1 \times 1} \times \frac{1}{4} \times \frac{1}{2} = \underline{\underline{\frac{3}{8}}} (= 0.375)$$

$$\therefore ③ \text{Mean} = np = 3 \times 0.5 = 1.5 \quad \text{Ans}$$

$$\therefore ④ \text{Variance} = npq = 3 \times 0.5 \times 0.5 = 0.75$$

$$\therefore \text{St. Dev} = \sqrt{npq} = \sqrt{0.75} = 0.8660$$

~~DISCRETE~~

- Poisson Distribution
- The Poisson Distribution is an discrete Prob. distribution that deals with the # Events that occurs within a fixed Interval of Time / Space given an Known "Average" Rate of Occurrences.
- No. of Events occurring in a fixed Time Interval.

$$P(X) = \frac{\lambda^x e^{-\lambda}}{x!}$$

- # Imp. Points about Poisson Distn
- Poisson Distn is an pmf i.e.  $\sum P(X) = 1$ .
  - Moment Generating Fn (MGF) =  $e^{\lambda(e^t - 1)}$
  - Characteristic Fn (Chf Fn) =  $e^{\lambda(e^{it} - 1)}$
  - Prob. Generating Fn (PGF),  $G(t) = e^{\lambda(t-1)}$
  - In Poisson Distn, Mean = Variance =  $\lambda t$  ( $t = \# \text{Intervals}$ )

- Q → The Avg. Number of Customers Entering a Store In an Hour is 5. What is the Prob. That exactly 3 customers will enter the next hour?

→ Here,  $\lambda = 5$

$$\therefore P(X=3) = \frac{\lambda^x e^{-\lambda}}{x!} = \frac{5^3 (2.7182)^{-5}}{3!} = 0.142$$

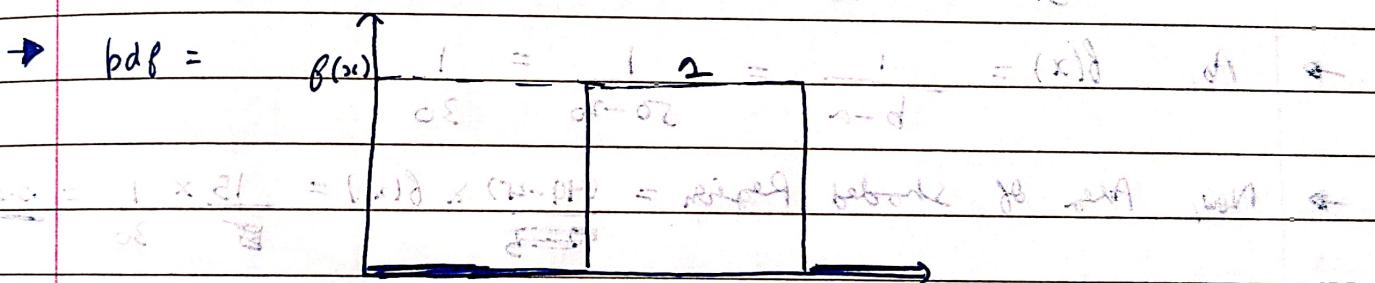
Q III ly, Mean =  $\lambda t = 5$  ( $t = 1$ )  $\rightarrow 5$   
Variance =  $\lambda t = \lambda (t=1) = 5$

## Continuous Uniform Distribution $\rightarrow$ CONTINUOUS...

- An Uniform Distribution refers to an Type of Prob. Distribution in which outcomes are been Equally Likely.
- Uniform Dist<sup>n</sup>  $\rightarrow$ 
  - Discrete Uniform Dist<sup>n</sup> (pmf).  $\rightarrow$  (See 3 Page Back)
  - Continuous Uniform Dist<sup>n</sup> (pdf).

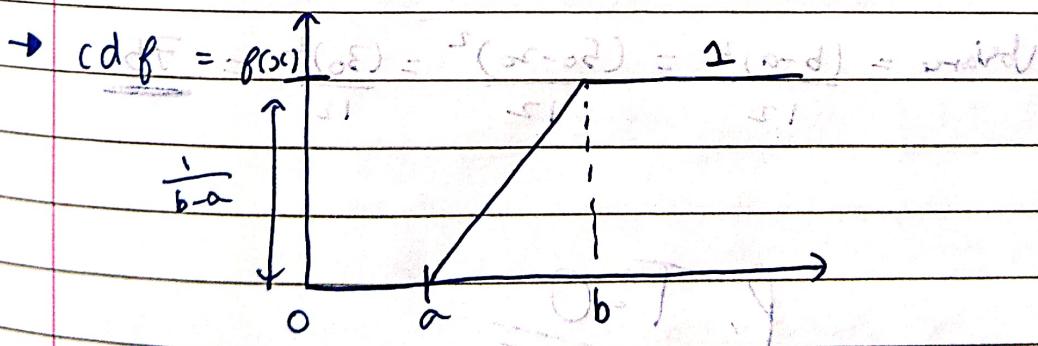
- An Continuous Distribution F<sup>n</sup> Is an dist<sup>n</sup> that Has an Infinite No. of Values defined In an Specified Range / Bound.
- The Random Variable Is Continuous In This.
- Notation :  $U(a, b)$

Parameters :  $a < b < \infty$ ,  $b > a$



$\rightarrow$  Area of Rectangle =  $L \times B = (b-a) \times f(x)$

$$1 = (b-a) \times f(x) \Rightarrow f(x) = \frac{1}{b-a}$$



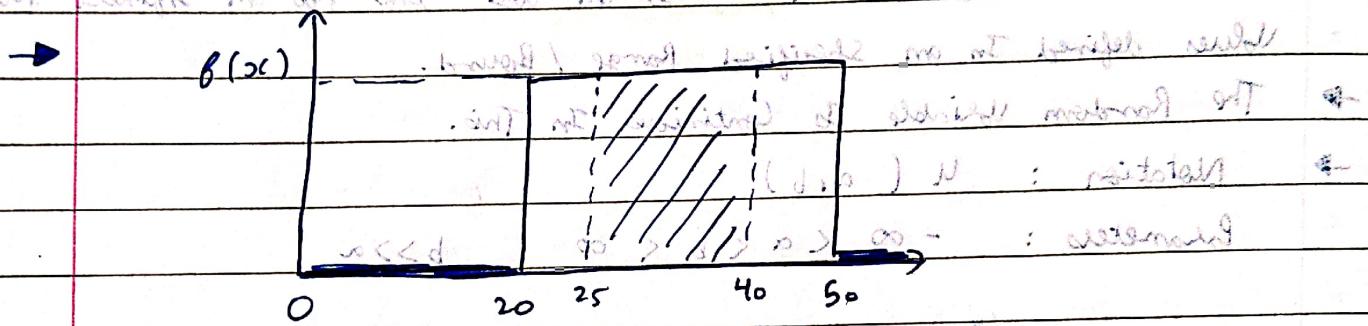
$\therefore$  Thus,  $CDF = \begin{cases} 0 & , x < a \\ \frac{x-a}{b-a} & , x \in [a, b] \\ 1 & , x > b \end{cases}$

Note → For Continuous Uniform Dist<sup>n</sup>, Mean =  $\frac{a+b}{2}$  } No Need for Loop  
 ... ~~PROGRAMMING~~ ← ~~initialising variables~~ } for AIML...

$$\text{Variance} = \frac{(b-a)^2}{12}$$

Q. → The No. of Items sold at a shop daily is uniformly distributed with the Max & Min Items sold as 50 & 20.

→ Find the Prob. of sales to fall b/w 25 & 40



$$\rightarrow \text{As, } f(x) = \frac{1}{b-a} = \frac{1}{50-20} = \frac{1}{30} \text{ l.h.s.} = \text{R.H.S.}$$

$$\rightarrow \text{Now, Area of shaded Region} = \frac{(40-25)}{30} \times f(x) = \frac{15}{30} \times \frac{1}{30} = \underline{\underline{0.5}}$$

∴ Thus, 50% chance that No. of Items sold in [25, 30]

$$\textcircled{1} \quad \text{For Range } a=20 \text{ and } b=50, (a+b) = 1 = \text{Mean}$$

$$\text{Mean} = \frac{b+a}{2} = \frac{50+20}{2} = \frac{70}{2} = \underline{\underline{35}}$$

$$\text{Variance} = \frac{(b-a)^2}{12} = \frac{(50-20)^2}{12} = \frac{(30)^2}{12} = \underline{\underline{75}}$$

CONTINUOUS

→ Normal Distribution / Gaussian Distr'n

→ An Continuous Prob. distribution. ( $\mu = 3 \times 3 = 9$ )

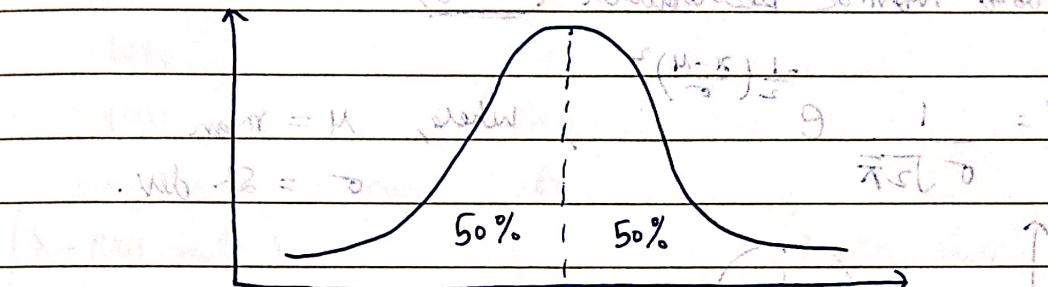
→ It is an Bell-shaped distribution:  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

→ Most of The Real world data follows on Normal distribution.

200-2A

Q) Eg:- Height of an Population, Exam Score of an class

(AN 2) ~~identical~~ somewhat identical



# Characteristics of Normal Dist.

→ ND is symmetric about the mean. Plain Normal.

$\rightarrow \text{Mean} = \text{Median} = \text{Mode}$ . Will become no real life as no ST + E

→ Skewness of ND is 0.

Left side of hill is a mix of sand & shale

## # Empirical Rule of a Normal Distribution

$$68\% = 95\% = 99.7\% \text{ Rule (sigma)}$$

(mean)  $68\%$   $\underline{\hspace{1cm}}$   $95\%$   $\underline{\hspace{1cm}}$   $99.7\%$  Rule  $\underline{\hspace{1cm}}$  (standard)

With regard to our own (I mean the training) I am still

68% of The Total data  99.7% of points lie within

points will lie within 95% of the total data 3σ Range of ND

1 st. dev. Range / points will lie within  $2\sigma$   $(x - 3\sigma, x + 3\sigma)$

$$68\% \text{ In } [\bar{x} - \sigma, \bar{x} + \sigma] \quad 95\% \text{ In } [\bar{x} - 2\sigma, \bar{x} + 2\sigma]$$

∴ Thus, Based on The Empirical Rule:-

$$P(\mu - \sigma \leq X \leq \mu + \sigma) = 68\% = 0.68$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 95\% = 0.95$$

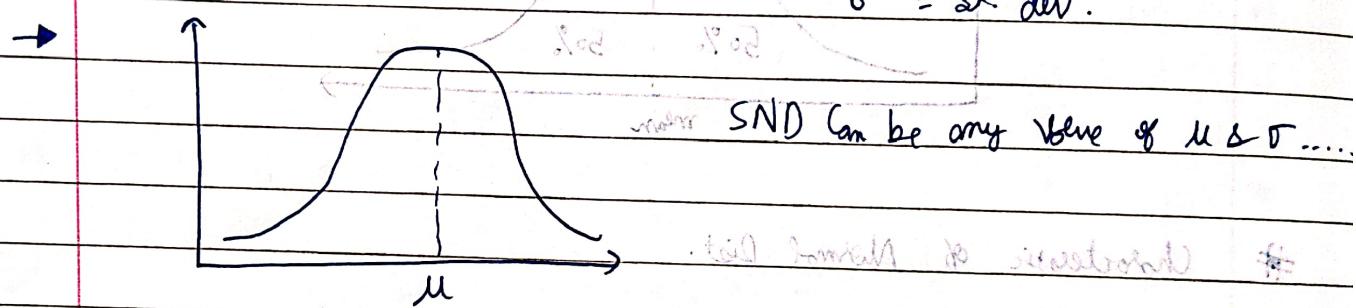
$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) = 99.7\% = 0.997$$

- JAI LOGIC

⇒ Standard Normal Distribution (SND)

$$\rightarrow P.d.f = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2}$$

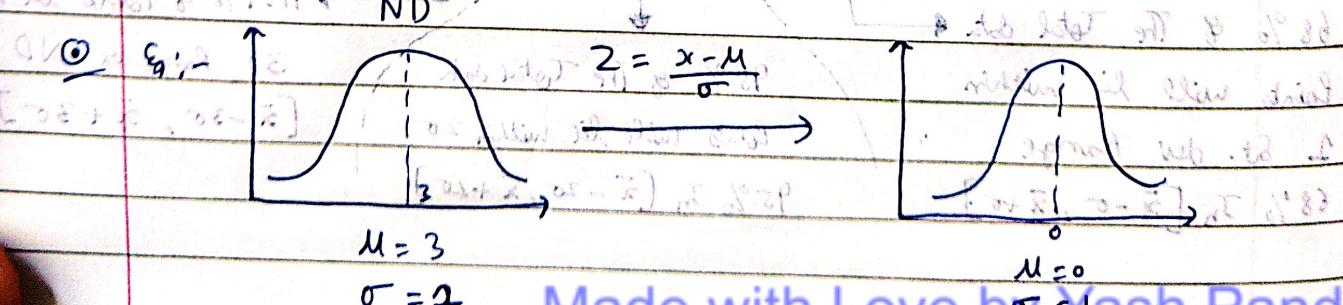
where,  $\mu = \text{mean}$   
 $\sigma = \text{st. dev.}$



- St. Normal Distn Is a Continuous Probability distribution.
- It Is a Special Case of Normal Distn where  $\mu=0$  &  $\sigma=1$ .

Note → Let, Graph 1 Represent Normal Distn Curve for Marks In Science Stream.  
 Graph 2 " " " " " Arts Stream.

→ Ab, There are Many different factors In Both Streams, Thus, We Can't Compare The Graph directly, Thus, For Comparison,  
 We will Convert ND Curve To SND Curve, Then Compare Them.



Q. Standardize the below data of Normal Distn:-

Given That  $\mu = 1010 \Delta \sigma = 20$

$x$	$x - \mu$	$Z = (x - \mu) / \sigma$
950	-60	-3
970	-40	-2
990	-20	-1
1010	0	0
1030	20	1

Normal distribution (050) and New N(1010) have mean 2 standard deviation

(original) center 1070 and distance 200 from 1010 to 1050 is 3 times the std dev

$(X\text{-Axis Scale})$  of ND  $\rightarrow$   $(X\text{-Axis Scale})$  of SND

Note  $\rightarrow Z = \frac{x - \mu}{\sigma}$  The Z-value simply represent that how much st.dev with mean distance at which the point is away from mean....

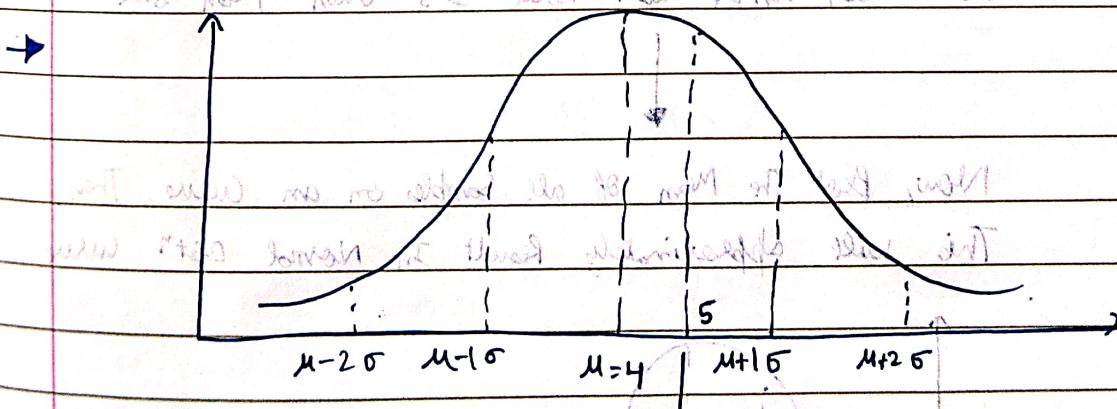
$\sigma = 0.5 \Delta z = 0.5 \Delta \mu$  point will remain

Q: -  $x = 5, \mu = 4, \sigma = 2$

Then,  $Z = \frac{x - \mu}{\sigma} = \frac{5 - 4}{2} = 0.5 \therefore$  The point is  $(0.5\sigma)$  away from mean

(standardized point) therefore 0.5 standard deviation off from the mean of ND....

After point was zed and then plotted on scale

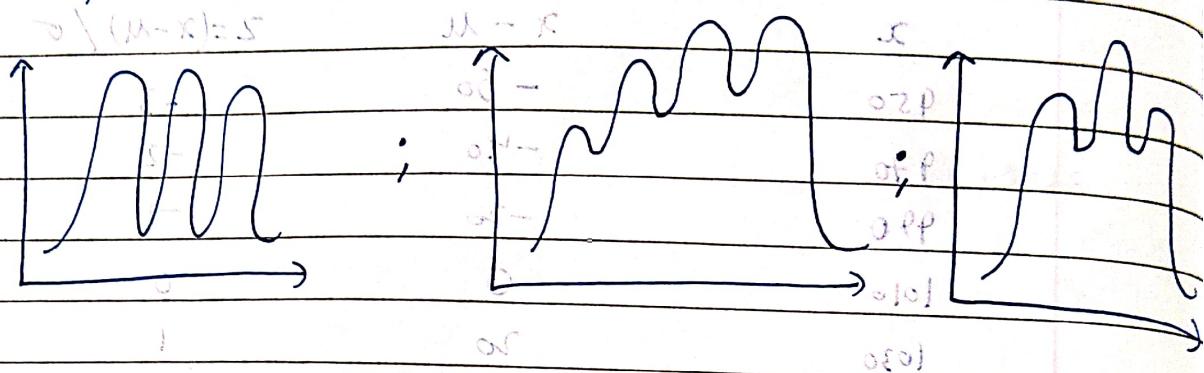


$\mu + 0.5\sigma : (0.5\sigma)$  away from mean

- JAI LOGIC

→ Central Limit Theorem (CLT)

→ Sometimes, The distributions can be Irregular.



→ The Central Limit Theorem states that If you have an population with a mean  $\mu$  and st. dev  $\sigma$  and take sufficiently large samples (random) from the population with replacement, Then, the distribution of the sample will be Approximately Normal Distn.

→ Sampling Mean of Population ( $\bar{x}$ ) will be Approx. To Normal distn.



Given, an population curve which is Irregular and Not

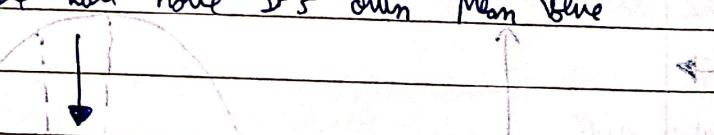
Normal Distn having Mean =  $\mu$  & St. dev =  $\sigma$

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

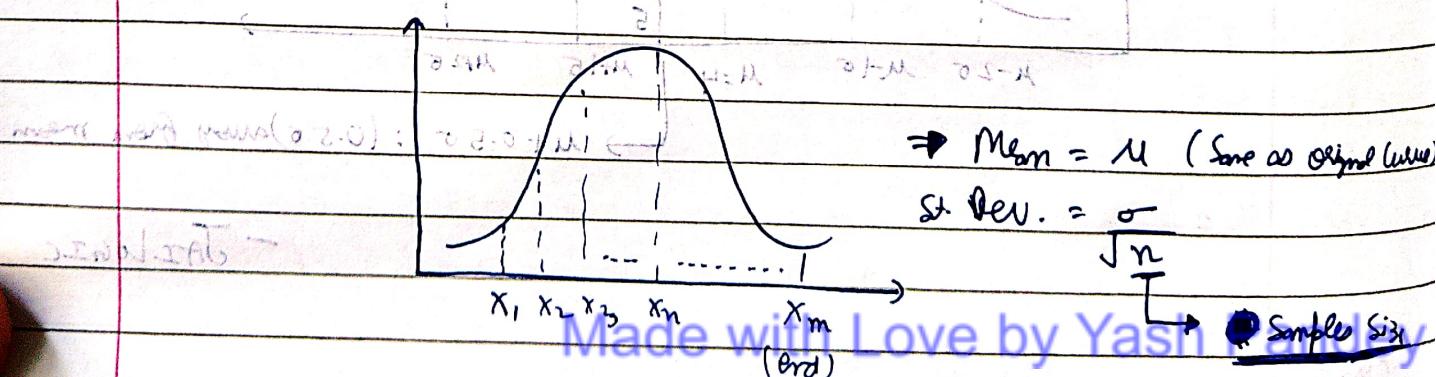
$$\text{and mean}(\bar{x}) = \mu \text{ and } \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Divide The Population Into Samples (large No. of Samples)

where Each Sample will have It's own Mean Value



Now, Plot The Mean of all Samples on an Curve, Then,  
This will approximately Result In Normal Distn Curve



## # The Two Conditions of CLT

- The # samples should be large.
- The Sample Size should be Greater Than Equal To 30.  
(Except The Population dist<sup>n</sup> which is already an Normal dist<sup>n</sup>)

Note → In Mean Curve, Mean =  $\mu$

$$\text{St. dev} = \frac{\sigma}{\sqrt{n}} \quad (\text{Standard Error})$$

Where, Standard Error of  $\frac{1}{\sqrt{n}}$   $\Rightarrow n \uparrow \text{SE} \downarrow$   
 $n \downarrow \text{SE} \uparrow$

Q → You have an population dist<sup>n</sup> with  $\mu = 100$  and  $\sigma = 20$ . If You Take sample size 50 from The population, Then, what will be Prob. That mean will be less Than 105?

$$\rightarrow \mu = 100, \sigma = 20 \text{ and } n = 50$$

→ ATQ, we have to Find Prob. Such That Sample Mean  $< 105$ .

$$\therefore \bar{X} = 105, \mu_s = \mu = 100$$

$$\sigma_s = \frac{\sigma}{\sqrt{n}} = \frac{20}{\sqrt{50}} = \frac{20}{5\sqrt{2}} = 2\sqrt{2} = 2.83$$

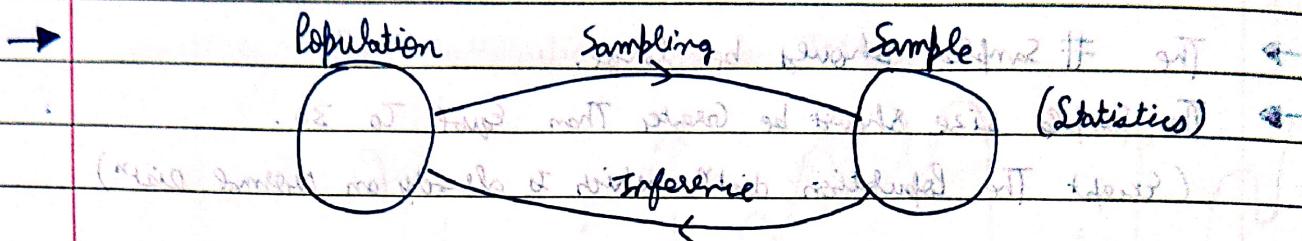
$$\therefore Z \text{ score} = \frac{\bar{X} - \mu_s}{\sigma_s} = \frac{105 - 100}{2\sqrt{2}} = \frac{5}{2\sqrt{2}} = 1.7675$$

(for Sample)

$$\therefore P(Z < 1.7675) = 0.8944 \quad (\text{Value from Z-Table})$$

$\therefore$  Thus, The probability will be 0.89.

## → Estimate In Statistics



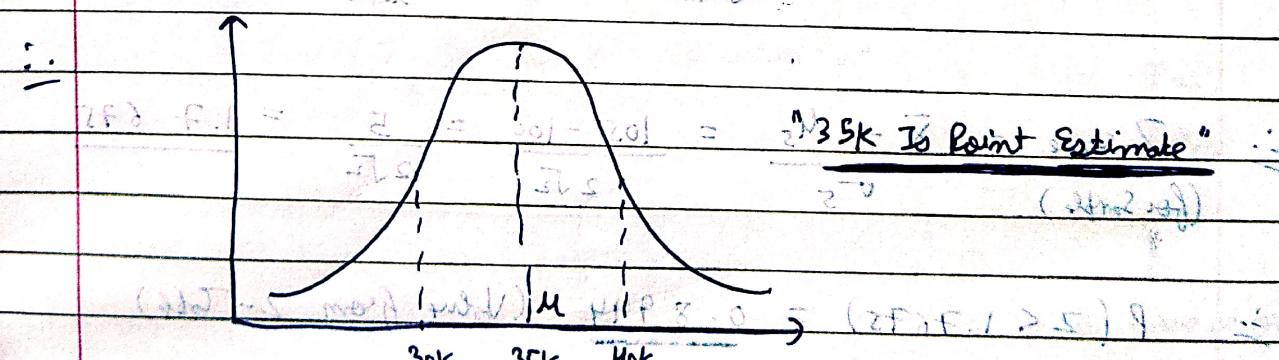
→ An Specified Observed Value which Is been Used To Estimate an Unknown Population Parameter Using The Samples.

- ① Point Estimate → An Point Estimate is a Single Value that Is been Used To Estimate The True Value of The Population Parameter.
- ② Eg:- Avg. Salary of IT-Employee of TCS.

- ③ Interval Estimate → Range of The Values used To Estimate The Unknown Population Parameters.

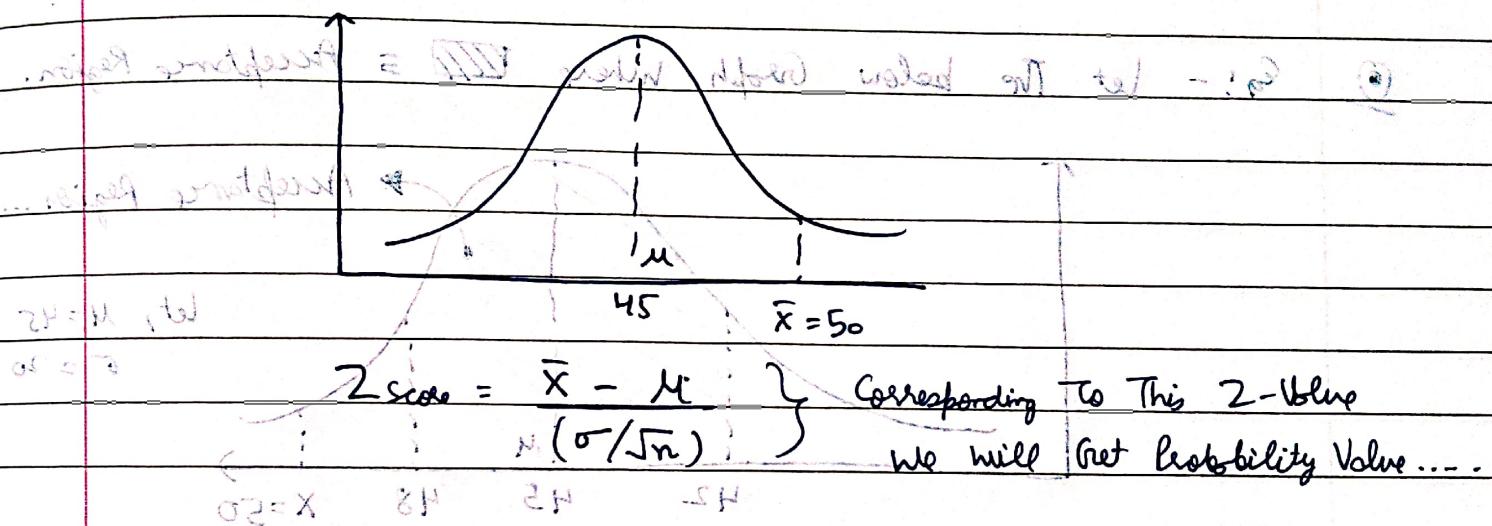
- ④ Eg:- The Avg. Salary Is 35K Whereas The Range (Confidence Range) where The Max Employee Lie Is 30k To 40k.

∴ Thus,  $\bar{x} = 35K$  Is The Point Estimate  
 $\bar{x} \pm 2\sigma = 35K \pm 5K = (30K, 40K)$  Is The Interval Estimate.



- P-Value ~~is nothing~~ is the probability of getting a test statistic as extreme as or more extreme than the observed value under the null hypothesis.
- The P-Value Is The Probability Value Calculated from an Statistical Test.
- P-Value Is Hypothesis Testing which is used to decide whether To Reject The Null Hypothesis or Not.

Q) Given - Mean of Employee Is 45 Years. Let,  $\bar{X}$  of Sample Comes So -



Q) Ex:- let, am Medicine Is 95% Effective  
which means, on 100 People, 95 will See Result whereas 5 won't

$$\text{Margin of Error} = 2H = \sigma = \mu - x̄ = 2(2.5)$$

∴ Thus, We Are 95% Confident & 5% Margin of Error.  
 $\therefore \text{Margin of Error} = (2(2.5)) = 5$

\* Level of Significance ( $\alpha$ ) = Margin of Error = 5 %

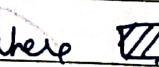
$$5 = \sigma \text{ which is } \Rightarrow \sigma = (2(2.5)) = 5$$

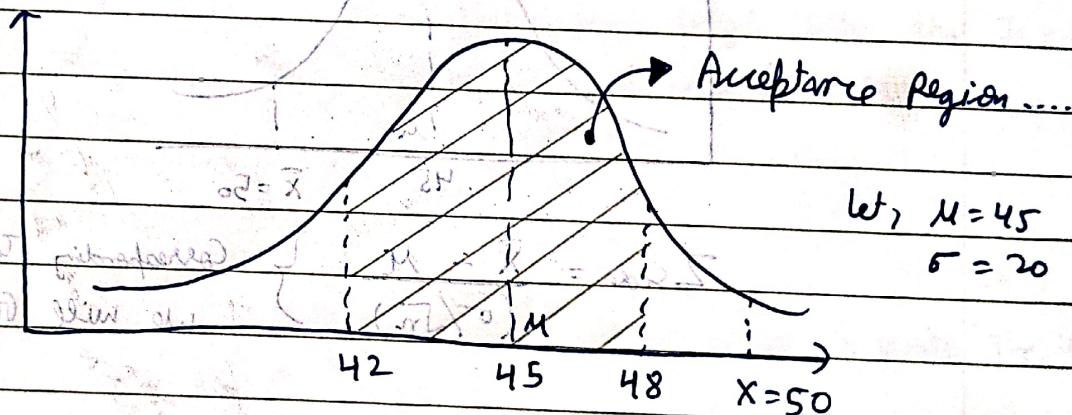
∴ Thus, Confidence Interval =  $1 - \alpha$  (Error)  $= 1 - 0.05 = 0.95$

Note → If the Value Comes out of The Confidence Value Range, Then, Calculate Z Value of That Value & Corresponding Prob. To That Z-Value.

→ If, (Prob.  $<$  level of sig.)  $\rightarrow$  Reject The Value  
 (Prob.  $\geq$  level of sig.)  $\rightarrow$  Accept. The Value

JAZ Logic

Q:- Let The below Graph where  = Acceptance Region.



→ Let, level of Sig ( $\alpha$ )  $= 2.5\% = 0.05$

Then, for The Value  $x=50$ , Calculate Z-score,

$$Z \text{ score} = \frac{x - \mu}{\sigma} = \frac{50 - 45}{20/\sqrt{n}} = \text{Something}$$

Then,  $P(Z \text{ score}) = \text{min}$

→ If, Prob (Z score)  $< \alpha \rightarrow$  Reject The Value  $x=50$

If, Prob (Z score)  $\geq \alpha \rightarrow$  Accept The Value  $x=50$

## ⇒ Applications of Z-Score

→ Let, for a Normal Distribution ( $\mu = 50, \sigma = 20$ ,  $X = 110$ )  
 Then, Let  $Z = \frac{X - \mu}{\sigma} = \frac{110 - 50}{20} = 3$  →  $X = 110$  is 3σ Away  
 A.J.F.  $Z = 3.27$  is approximately 3.5 σ away from the Mean....

① Let, we have to find  $P(Z > 0.5)$

$$\text{Then, } P(Z > 0.5) = 1 - P(Z \leq 0.5)$$

(approximate standard deviation)  $= 1 - 0.6915$  (from Z-table)

$$= 0.3185$$

② What is the  $P(\text{Marks less than } 3)$  when  $\mu = 40$  &  $\sigma = 1.5$

$$\rightarrow Z = \frac{(X - \mu)}{\sigma} = \frac{3 - 40}{1.5} = -1$$

$$\text{Then, } P(Z = -1) = 0.1587$$

∴ The  $P(\text{Marks less than } 3)$  is 0.1587.

$$\text{③ } \text{If } P(Z \leq 2.5) = P(Z < 3) = -P(Z > 2) \text{ (from CH)}$$

$$\therefore P(2 < Z \leq 3) = P(Z < 3) - P(Z < 2)$$

$$\text{f(H)} \quad f(z) = 0.9987 - 0.9772$$

$$\text{(H)} \quad = 0.0215 = 2.15\%$$

→ Hypothesis and Mechanism

→ Hypothesis Is an Claim, as a Statement / Assumption about The Population Parameters That Can be Tested Using The Statistical Methods.

Ex:- Avg Salary of Employee of TCS Is 7 LPA.

① Null Hypothesis : The Initial / default Assumption.

② Alternate Hypothesis : The Opposite of The Null Hypothesis.

# Mechanism of Hypothesis (How it work about) :-

1. Form The Hypothesis

→ Null Hypothesis ( $H_0$ ) :-  $\mu = 5$

→ Alternate Hypothesis ( $H_A$ )

Eg:- The Age of PW Skills Is Atleast 35 Years.

$H_0$  will always have = (Sign) 19       $\mu \geq 35$        $H_A$

Eg:- The Avg Salary of TCS Employee Is 7 LPA

$\mu = 7$

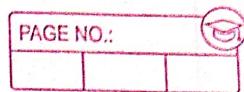
$\mu \neq 7$

$\mu = 7$        $H_0$

$H_A$

Eg:- The Age for M.Tech Is Almost 30 Years.

$\mu \leq 30$        $\mu > 30$



② Experiment / Statistical Analysis. ( p-Value / Significance Value)

Test statistic value and p-value both give the probability which are used for decision.

③ Reject H<sub>0</sub> / Fail to Reject H<sub>0</sub>. i.e. decide with p-value

p-value (p-value <=) is called robustness of the test distribution and significance level

Decision point: If the p-value is less than the significance level then reject the null hypothesis.

$\Rightarrow$  Hypothesis Testing vs Statistical Analysis

(Z-test =  $\frac{x - \mu}{\sigma}$ ) compared to the level of significance

→ Z-Test.  $\rightarrow$  Average ...

→ T-Test. (H<sub>0</sub>:  $\mu_1 = \mu_2$ )  $\rightarrow$  H<sub>1</sub>:  $\mu_1 \neq \mu_2$  based on sample size

→ Chi-Square Test  $\rightarrow$  Categorise Data ...

→ ANNOVA.  $\rightarrow$  Variance ...

→ A/C To Central Limit Theorem,  $\sigma = (\infty) \cdot 0.12$  at level 0.05

→ Z-Test ( $H_0: \mu \leq \mu_0$ ) & ( $H_1: \mu > \mu_0$ )  $\rightarrow$  p-value

→ A/C To Central Limit Theorem,  $\sigma = (\infty) \cdot 0.12$  at level 0.05

$\mu_s = \mu$  and  $\sigma_s = \frac{\sigma}{\sqrt{n}}$

$\{ \text{if } n \geq 30 \} \rightarrow Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \rightarrow$  p-value

→ For Sample Z-Test, Sample size  $\geq 30$

$\sigma_p$  must be known

$$\text{Then, } Z \text{ score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = 20.25$$

∴ Then, 5 steps are

- Form The Hypothesis.
- Level of Significance.
- Type of Test Needed.
- Z Score Evaluation.
- Probability

(Contd...)

- Q → Suppose an Child Physiologist says That The Avg. Time Working mother Spends Talking Their Children Is ~~11.5~~ upto 11 minutes per day. To Test The Hypothesis, We Conducted The Exp. on Random Sample of 100 Working mothers and found That They Spend 11.5 min per day with Their Children. Assume That The Pop.  $\sigma = 2.3$  min. Then, Conduct The Test with The 5 % level of significance ( $\alpha = 0.05$ ).

→ Time upto 11 min  $\rightarrow H_0: \mu \leq 11$

$H_A: \mu > 11$

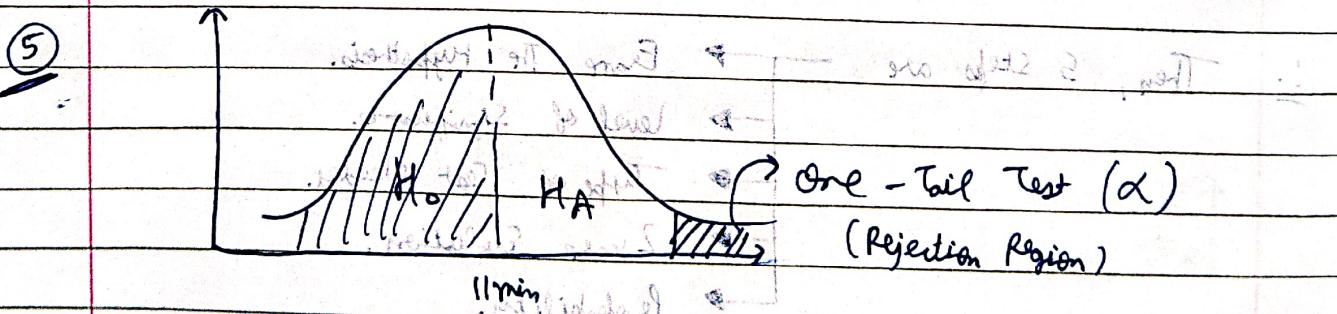
① Form The Hypothesis

→ Then,  $H_0: \mu \leq 11$  &  $H_A: \mu > 11$

② Level of Sig. ( $\alpha$ ) = 0.05

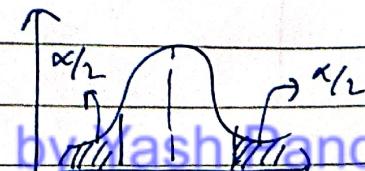
③ Type of Test  $\rightarrow$  Z-Test [As,  $n > 30$ ]  
Given  $\sigma = 2.3$

$$Z_{\text{Score}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{11.5 - 11}{2.3/\sqrt{100}} = 2.17$$



Note → Whenever,  $H_A$  has  $\mu \neq \mu_0$ , Then,

It will be Two Tail Test





→ Now for checking of Prob. 1 we have 2 methods:-

PREFER  
(i)

$Z_{\text{critical}} \approx \alpha = 0.05$  (from Z-Table, find The closest Z-value)  
To This, Then, Take Z-value of Test

$$Z_{\text{critical}} = 1.64 \quad \rightarrow \quad Z \leq 1.64 \quad (\text{Accept})$$

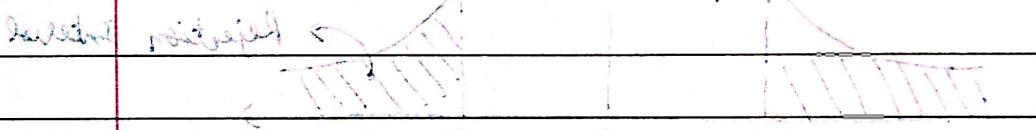
$$\underline{18.5} = \underline{8.0 + 10.5} > 1.64 \quad (\text{Reject})$$

→ As,  $Z_{\text{score}} = 2.17 > 1.64 \rightarrow \text{Reject The } H_0 \dots$

(ii) p-Value :  $P(Z_{\text{score}} = 2.17) = 0.9850$

$$\text{Actual significance} P(Z > 2.17) = 1 - 0.9850 = 0.015 < \alpha$$

As,  $0.015 < \alpha$ , Reject The  $H_0 \dots$



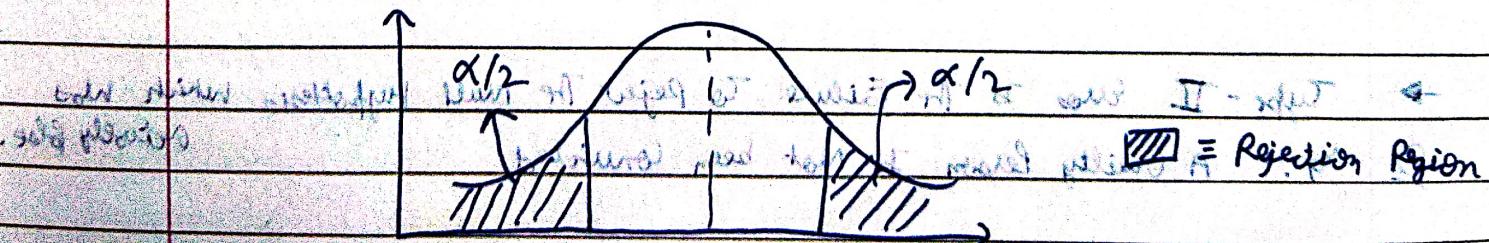
Q → The Average Height of all Residents In an city Is 168cm with an  $\sigma = 3.9$  cm. On Research, They Believe That mean To be different. He measures The Height of 36 Individuals and found as 169.5 cm. Test The Hypothesis with 95% Confidence Interval.

→ 95% Confidence Interval  $\rightarrow \alpha = 0.05$

(i) Frame The Hypothesis  $\rightarrow H_0: \mu = 168$  (H<sub>0</sub>)

$H_A: \mu \neq 168$  (H<sub>A</sub>)

→ As, H<sub>A</sub> has ( $\neq$ ) Sign, Thus, This will be an Two-Tail Test.

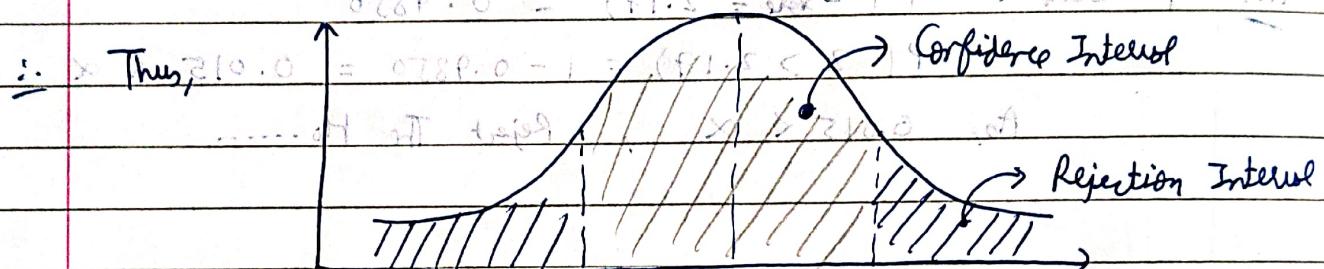


② Level of Significance : As, It is an Two-Tail Test,  
 $\alpha = 2.5\% = 0.025$

③ Type of Test : Z - Test

$$\text{Z score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{169.5 - 168}{3.9/\sqrt{36}} = 2.31$$

⑤  $Z_{\text{critical}}$  : Here,  $\alpha = 0.025$ , Then,  
 $Z_{\text{critical}} = 1.96$  (Evaluated from Z-Table)



$Z_{\text{score}} = 2.31 > 1.96$  → Reject The  $H_0$ .

∴ Thus, we Reject The Null Hypothesis.

EASY

→ Type-I and Type-II Errors

→ Type-I Error Is The Rejection of The Null Hypothesis when It Is been actually True.

① Eg:- An Innocent Person Is been Convicted.

→ Type-II Error Is The Failure To Reject The Null Hypothesis which Was

② Eg:- A Guilty Person Is Not been Convicted.



## True Nature of Hypothesis (Actual)

	H <sub>0</sub> Is True	H <sub>0</sub> Is False	H <sub>A</sub> Is True	H <sub>A</sub> Is False
Conclusion (Based on The Analysis)	Support H <sub>0</sub> Support H <sub>A</sub>	Accept H <sub>0</sub> Reject H <sub>0</sub>	Correct Conclusion Type - I Error	Type - II Error Correct Conclusion

$$\alpha = \left( \text{Type - I Error} \right) \text{ PVA}$$

$\text{H}_0$  Is True  $\xrightarrow{\quad}$  Accept  $\text{H}_0$   $\xrightarrow{\quad}$  Correct Conclusion

$\text{H}_0$  Is True  $\xrightarrow{\quad}$  Reject  $\text{H}_0$   $\xrightarrow{\quad}$  Type - I Error

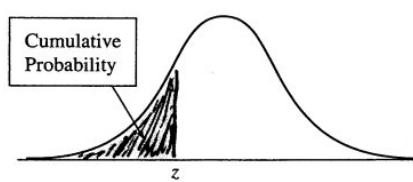
$\text{H}_0$  Is False  $\xrightarrow{\quad}$  Accept  $\text{H}_0$   $\xrightarrow{\quad}$  Type - II Error

$\text{H}_0$  Is False  $\xrightarrow{\quad}$  Reject  $\text{H}_0$   $\xrightarrow{\quad}$  Correct Conclusion

Hypothesis Test (H.T.) : Null - or  $H_0$  - JAI Logic

Hypothesis Test (H.T.)

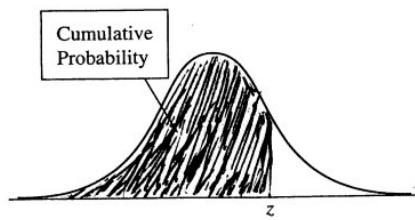
# APPENDIX A



Cumulative probability for  $z$  is the area under the standard normal curve to the left of  $z$

**TABLE A Standard Normal Cumulative Probabilities**

<b><i>z</i></b>	<b>.00</b>	<b><i>z</i></b>	<b>.00</b>	<b>.01</b>	<b>.02</b>	<b>.03</b>	<b>.04</b>	<b>.05</b>	<b>.06</b>	<b>.07</b>	<b>.08</b>	<b>.09</b>
-5.0	.000000287	-3.4	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0002
-4.5	.00000340	-3.3	.0005	.0005	.0005	.0004	.0004	.0004	.0004	.0004	.0004	.0003
-4.0	.0000317	-3.2	.0007	.0007	.0006	.0006	.0006	.0006	.0006	.0005	.0005	.0005
-3.5	.000233	-3.1	.0010	.0009	.0009	.0009	.0008	.0008	.0008	.0008	.0007	.0007
		-3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0011	.0010	.0010
		-2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
		-2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
		-2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
		-2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
		-2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
		-2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
		-2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
		-2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
		-2.1	.0179	.0174	.0170	.0166	.0162	.0158	.0154	.0150	.0146	.0143
		-2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
		-1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
		-1.8	.0359	.0351	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
		-1.7	.0446	.0436	.0427	.0418	.0409	.0401	.0392	.0384	.0375	.0367
		-1.6	.0548	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
		-1.5	.0668	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
		-1.4	.0808	.0793	.0778	.0764	.0749	.0735	.0721	.0708	.0694	.0681
		-1.3	.0968	.0951	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
		-1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
		-1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
		-1.0	.1587	.1562	.1539	.1515	.1492	.1469	.1446	.1423	.1401	.1379
		-0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611
		-0.8	.2119	.2090	.2061	.2033	.2005	.1977	.1949	.1922	.1894	.1867
		-0.7	.2420	.2389	.2358	.2327	.2296	.2266	.2236	.2206	.2177	.2148
		-0.6	.2743	.2709	.2676	.2643	.2611	.2578	.2546	.2514	.2483	.2451
		-0.5	.3085	.3050	.3015	.2981	.2946	.2912	.2877	.2843	.2810	.2776
		-0.4	.3446	.3409	.3372	.3336	.3300	.3264	.3228	.3192	.3156	.3121
		-0.3	.3821	.3783	.3745	.3707	.3669	.3632	.3594	.3557	.3520	.3483
		-0.2	.4207	.4168	.4129	.4090	.4052	.4013	.3974	.3936	.3897	.3859
		-0.1	.4602	.4562	.4522	.4483	.4443	.4404	.4364	.4325	.4286	.4247
		-0.0	.5000	.4960	.4920	.4880	.4840	.4801	.4761	.4721	.4681	.4641



Cumulative probability for  $z$  is the area under the standard normal curve to the left of  $z$

**TABLE A Standard Normal Cumulative Probabilities (continued)**

<i>z</i>	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
0.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
0.2	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141
0.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
0.4	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879
0.5	.6915	.6950	.6985	.7019	.7054	.7088	.7123	.7157	.7190	.7224
0.6	.7257	.7291	.7324	.7357	.7389	.7422	.7454	.7486	.7517	.7549
0.7	.7580	.7611	.7642	.7673	.7704	.7734	.7764	.7794	.7823	.7852
0.8	.7881	.7910	.7939	.7967	.7995	.8023	.8051	.8078	.8106	.8133
0.9	.8159	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.8461	.8485	.8508	.8531	.8554	.8577	.8599	.8621
1.1	.8643	.8665	.8686	.8708	.8729	.8749	.8770	.8790	.8810	.8830
1.2	.8849	.8869	.8888	.8907	.8925	.8944	.8962	.8980	.8997	.9015
1.3	.9032	.9049	.9066	.9082	.9099	.9115	.9131	.9147	.9162	.9177
1.4	.9192	.9207	.9222	.9236	.9251	.9265	.9279	.9292	.9306	.9319
1.5	.9332	.9345	.9357	.9370	.9382	.9394	.9406	.9418	.9429	.9441
1.6	.9452	.9463	.9474	.9484	.9495	.9505	.9515	.9525	.9535	.9545
1.7	.9554	.9564	.9573	.9582	.9591	.9599	.9608	.9616	.9625	.9633
1.8	.9641	.9649	.9656	.9664	.9671	.9678	.9686	.9693	.9699	.9706
1.9	.9713	.9719	.9726	.9732	.9738	.9744	.9750	.9756	.9761	.9767
2.0	.9772	.9778	.9783	.9788	.9793	.9798	.9803	.9808	.9812	.9817
2.1	.9821	.9826	.9830	.9834	.9838	.9842	.9846	.9850	.9854	.9857
2.2	.9861	.9864	.9868	.9871	.9875	.9878	.9881	.9884	.9887	.9890
2.3	.9893	.9896	.9898	.9901	.9904	.9906	.9909	.9911	.9913	.9916
2.4	.9918	.9920	.9922	.9925	.9927	.9929	.9931	.9932	.9934	.9936
2.5	.9938	.9940	.9941	.9943	.9945	.9946	.9948	.9949	.9951	.9952
2.6	.9953	.9955	.9956	.9957	.9959	.9960	.9961	.9962	.9963	.9964
2.7	.9965	.9966	.9967	.9968	.9969	.9970	.9971	.9972	.9973	.9974
2.8	.9974	.9975	.9976	.9977	.9977	.9978	.9979	.9979	.9980	.9981
2.9	.9981	.9982	.9982	.9983	.9984	.9984	.9985	.9985	.9986	.9986
3.0	.9987	.9987	.9987	.9988	.9988	.9989	.9989	.9989	.9990	.9990
3.5	.999767									
4.0	.9999683									
4.5	.9999966									
5.0	.999999713									

(standardized also called  $t$ )

## → Student - t - Distribution ( $t$ - dist<sup>n</sup>)

( $t$ -dist) used when  $\sigma_{\text{pop}}$  is not given & sample size  $n \leq 30$ .

→ So, for  $Z$ -Test,  $\sigma_{\text{pop}}$  must be given and  $n \geq 30$  is required.

→ But, In majority of The Cases ( $\sigma_{\text{pop}}$ ) won't be known.

∴ whenever,  $n \leq 30$  we use  $t$ -Test ( $t$  - dist<sup>n</sup>)

But, as  $t$ -Test st.  $\sigma_{\text{pop}}$  not given, it's related with  $s$  (sample std. dev.)

∴  $t$ -Test is not based on 'normal' method as defined normally we do  $Z$ -Test

→ for  $Z$ -Test,  $Z$  score =  $\frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$ ,  $\sigma$  not given

→ for  $T$ -Test,  $T$  -  $T$  score =  $\frac{\bar{x} - \mu}{s/\sqrt{n}}$  or  $\geq 2$  sig. diff.

$s/\sqrt{n}$  (min. val. of  $s^2$ )

Sample's St. dev...

Note → Degree of freedom of T-Test Is  $n-1$ .

→ Let, There are 5 Values, we have to select any 5 values such that the Avg. of These 5 Values is 10.

$$\text{Avg. } (\underline{\hspace{2cm}} \underline{\hspace{2cm}} \underline{\hspace{2cm}} \underline{\hspace{2cm}}) = 10$$

We can take any values for these 4 spots of  $\underline{\hspace{2cm}}$  if we have to take this value correctly and define our choice then we have to take this value such that Avg will be 10.

∴ Thus, for  $n$ -Values,  $(n-1)$  Values are Unconstrained  
2 Values are Constrained

∴ Thus, The degree of freedom will be  $(n-1)$ .

Q → Suppose a chief physiologist says that the average time working mothers spend talking their children is upto 11 min. To test the hypothesis, you took a random sample 20 working mothers and found that the avg. time spent is 11.5. The Sample St. Dev. is  $2.3$  ( $\alpha = 0.05$ )

→ Sample Size  $< 30$  Thus, Apply T-Test  
 $\sigma_{\text{pop}}$  is Not Given

① Frame The Hypothesis

$$\rightarrow M \leq 11 \rightarrow H_0 = M \leq 11 \\ H_A = M > 11$$

(No ≠ sign in  $H_A$ , Thus, One-Tailed Test)

② Level of Sig. ( $\alpha$ ) ( $= 0.05$ ) (One-Tailed Test)

③ Type of Test: T-Test

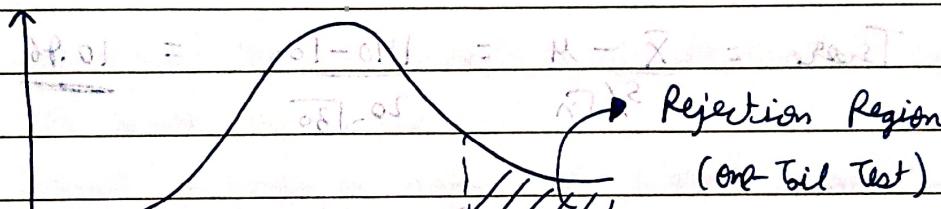
$$T\text{Score} = \frac{\bar{X} - M}{S/\sqrt{n}} = \frac{11.5 - 11}{2.3/\sqrt{20}} = 0.97$$

⑤ Test Score: Ans,  $\alpha = 0.05$

for This Question, degree of freedom =  $n - 1 = 20 - 1 = 19$

→ From T-Table with  $\alpha = 0.05$  & dof = 19, Tscore = 1.729

T  $\geq$  1.729  $\rightarrow$  Reject H<sub>0</sub>



→ Ans, Tscore = 0.97  $<$  1.729  $\rightarrow$  Accept The H<sub>0</sub>.....

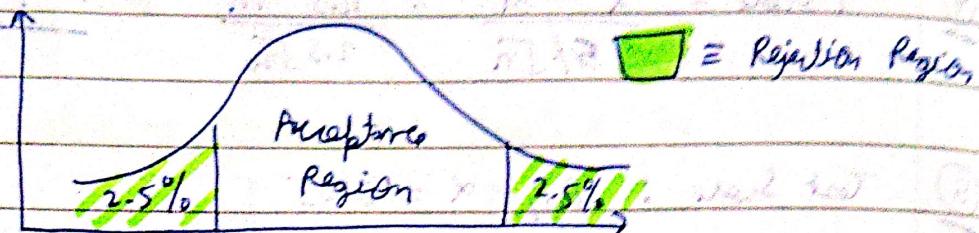
This, we failed to reject The Null Hypothesis.

→ In The Population, The Avg. IQ Is 100. A Team of Researchers Want To Test on Medicine To Check The +ve & -ve Effect on IQ. A Sample of 30 Participants Who Took Medicine Has Mean Value of IQ as 110 with St.dev. of 20. Did The medicine affect The IQ? ( $\alpha = 0.05$ )

→ Sample Size  $<$  30  
 $\sigma_{\text{Pop}}$  Is Not Given But  
 $\sigma_{\text{Sample}}$  Is Given

① Form The Hypothesis  $H_0: \mu = 100$  (H<sub>0</sub>)  $\rightarrow$   $H_1: \mu \neq 100$  (H<sub>1</sub>)  $\rightarrow$  Two-Tailed Test

② level of Sig. ( $\alpha$ ) = 0.05 (Two-Tailed Test)

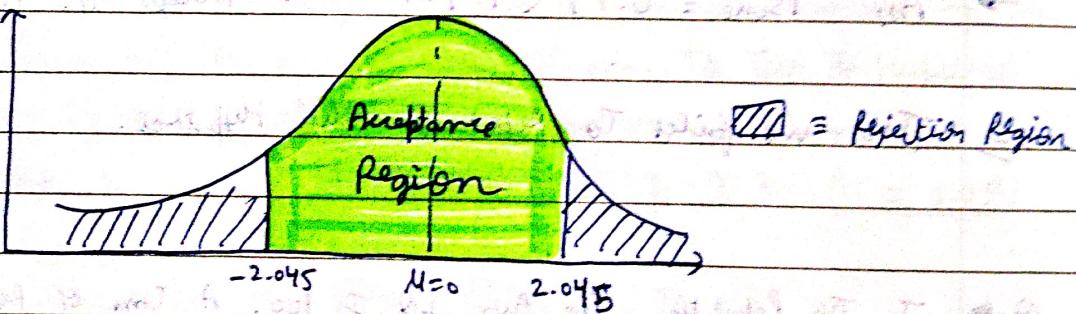


$\rightarrow$  degree of freedom (dof) =  $n - 1 = 29$

③ Type of Test : T-Test

$$\text{④ } T_{\text{Score}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{140 - 100}{5 / \sqrt{30}} = 10.96$$

⑤  $\alpha = 0.05$  (Two-Tail) } Then, from The T-Table  
 $dof = 29$  }  $T_{\text{Critical}} = 2.045$



$\rightarrow$  Now,  $T_{\text{Score}} = 10.96 > 2.045 \rightarrow$  Reject The H<sub>0</sub>....

$\therefore$  Thus, We Rejected The H<sub>0</sub>. I think we can't say it is 100% correct.

$\therefore$  Thus, The medicine has affected The IQ.

**t Table**

cum. prob	$t_{.50}$	$t_{.75}$	$t_{.80}$	$t_{.85}$	$t_{.90}$	$t_{.95}$	$t_{.975}$	$t_{.99}$	$t_{.995}$	$t_{.999}$	$t_{.9995}$
one-tail	<b>0.50</b>	<b>0.25</b>	<b>0.20</b>	<b>0.15</b>	<b>0.10</b>	<b>0.05</b>	<b>0.025</b>	<b>0.01</b>	<b>0.005</b>	<b>0.001</b>	<b>0.0005</b>
two-tails	<b>1.00</b>	<b>0.50</b>	<b>0.40</b>	<b>0.30</b>	<b>0.20</b>	<b>0.10</b>	<b>0.05</b>	<b>0.02</b>	<b>0.01</b>	<b>0.002</b>	<b>0.001</b>
df											
1	0.000	1.000	1.376	1.963	3.078	6.314	12.71	31.82	63.66	318.31	636.62
2	0.000	0.816	1.061	1.386	1.886	2.920	4.303	6.965	9.925	22.327	31.599
3	0.000	0.765	0.978	1.250	1.638	2.353	3.182	4.541	5.841	10.215	12.924
4	0.000	0.741	0.941	1.190	1.533	2.132	2.776	3.747	4.604	7.173	8.610
5	0.000	0.727	0.920	1.156	1.476	2.015	2.571	3.365	4.032	5.893	6.869
6	0.000	0.718	0.906	1.134	1.440	1.943	2.447	3.143	3.707	5.208	5.959
7	0.000	0.711	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	0.000	0.706	0.889	1.108	1.397	1.860	2.306	2.896	3.355	4.501	5.041
9	0.000	0.703	0.883	1.100	1.383	1.833	2.262	2.821	3.250	4.297	4.781
10	0.000	0.700	0.879	1.093	1.372	1.812	2.228	2.764	3.169	4.144	4.587
11	0.000	0.697	0.876	1.088	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	0.000	0.695	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.930	4.318
13	0.000	0.694	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.852	4.221
14	0.000	0.692	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.787	4.140
15	0.000	0.691	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733	4.073
16	0.000	0.690	0.865	1.071	1.337	1.746	2.120	2.583	2.921	3.686	4.015
17	0.000	0.689	0.863	1.069	1.333	1.740	2.110	2.567	2.898	3.646	3.965
18	0.000	0.688	0.862	1.067	1.330	1.734	2.101	2.552	2.878	3.610	3.922
19	0.000	0.688	0.861	1.066	1.328	1.729	2.093	2.539	2.861	3.579	3.883
20	0.000	0.687	0.860	1.064	1.325	1.725	2.086	2.528	2.845	3.552	3.850
21	0.000	0.686	0.859	1.063	1.323	1.721	2.080	2.518	2.831	3.527	3.819
22	0.000	0.686	0.858	1.061	1.321	1.717	2.074	2.508	2.819	3.505	3.792
23	0.000	0.685	0.858	1.060	1.319	1.714	2.069	2.500	2.807	3.485	3.768
24	0.000	0.685	0.857	1.059	1.318	1.711	2.064	2.492	2.797	3.467	3.745
25	0.000	0.684	0.856	1.058	1.316	1.708	2.060	2.485	2.787	3.450	3.725
26	0.000	0.684	0.856	1.058	1.315	1.706	2.056	2.479	2.779	3.435	3.707
27	0.000	0.684	0.855	1.057	1.314	1.703	2.052	2.473	2.771	3.421	3.690
28	0.000	0.683	0.855	1.056	1.313	1.701	2.048	2.467	2.763	3.408	3.674
29	0.000	0.683	0.854	1.055	1.311	1.699	2.045	2.462	2.756	3.396	3.659
30	0.000	0.683	0.854	1.055	1.310	1.697	2.042	2.457	2.750	3.385	3.646
40	0.000	0.681	0.851	1.050	1.303	1.684	2.021	2.423	2.704	3.307	3.551
60	0.000	0.679	0.848	1.045	1.296	1.671	2.000	2.390	2.660	3.232	3.460
80	0.000	0.678	0.846	1.043	1.292	1.664	1.990	2.374	2.639	3.195	3.416
100	0.000	0.677	0.845	1.042	1.290	1.660	1.984	2.364	2.626	3.174	3.390
1000	0.000	0.675	0.842	1.037	1.282	1.646	1.962	2.330	2.581	3.098	3.300
Z	0.000	0.674	0.842	1.036	1.282	1.645	1.960	2.326	2.576	3.090	3.291
	0%	50%	60%	70%	80%	90%	95%	98%	99%	99.8%	99.9%
	Confidence Level										

## # Z - Test Vs T - Test



If The St. dev. of Population ( $\sigma_{\text{Pop}}$ ) Is Given

Yes ↗ No ↘

Is The Sample Size  $\geq 30$

↗ No

Use T - Test

↙ Yes

Z - Test

- JAI LOGIC

Confidence Interval and Margin of Error

- An Estimate of pop. Population Parameter is an Approximation depending Solely on The Sample Information.
- A Point Interval Is Single No. whereas Conf. Interval Is an Interval.



Point Estimate Is The Exact Middle Location of The Confidence Interval.

Confidence Interval → 30 - 40

Point Estimate → 35



Confidence Interval (CI) = Point Estimate  $\pm$  Margin of Error



Z - Test Conf. Interval for Two - Tail Test

$$\rightarrow CI = \bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma_p}{\sqrt{n}} \rightarrow \text{Population St. Dev...}$$

↓

(Z Value at  $\alpha/2$ )

② Z-Test Confidence Interval for One-Tailed Test

$$\rightarrow CI = \bar{X} \pm Z_{\alpha} \cdot \frac{\sigma_p}{\sqrt{n}}$$

↓  
(Z-value at  $\alpha$ )

③ T-Test  $\rightarrow CI = \bar{X} \pm t_{\alpha/2} \cdot S$

EASY

Q → In An Exam, The St.dev. of marks is 100. A Sample of 36 students has an mean of 500. Construct The 95% of Confidence Interval about mean.

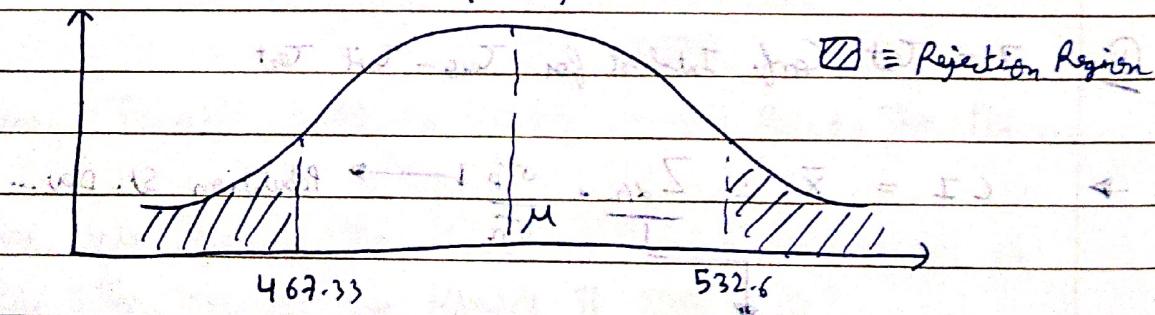
→ Level of Sig ( $\alpha$ ) =  $1 - (CI) = 1 - 0.95 = 0.05$

∴ Then,  $Z_{\alpha/2} = Z \left( \text{Value} = \frac{0.05}{2} \right) = 1.96$

∴  $CI = \bar{X} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$

① Upper CI =  $500 + 1.96 \cdot \left( \frac{100}{\sqrt{36}} \right) = 532.6$

② Lower CI =  $500 - 1.96 \left( \frac{100}{\sqrt{36}} \right) = 467.33$



∴ I am 95% Confident That Score In The Exam lies b/w [467.33 - 532.6].

→ Chi-Square dist<sup>n</sup> and Chi-Square Test

→ The Chi-Square dist<sup>n</sup> of an Probability Dist<sup>n</sup> that describes the sum of squares of K-Random Variables.

→ No Two-Tail Test Possible In Chi-Square Test.

$$H_0 = 11 - 2f = \text{constant} \quad H_1: f < 11$$

# The Properties of Chi-Square Dist<sup>n</sup>:

(a) Unimodal

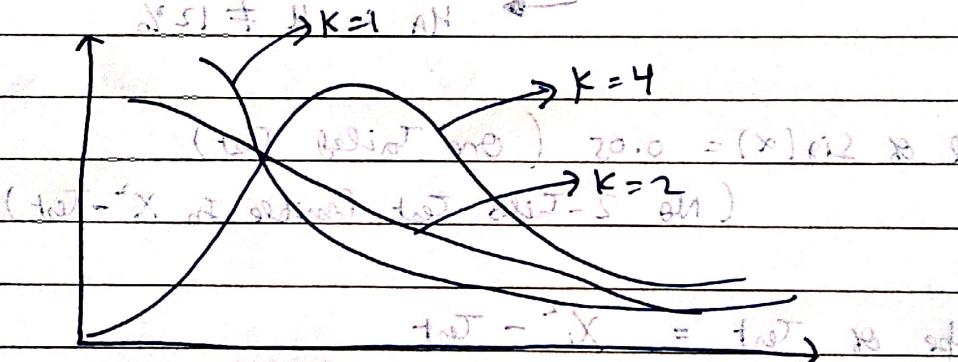
(b) Leptokurtic

→ The degrees of freedom (dof) is n-1.

→ Chi-Square dist<sup>n</sup> shape is determined by 'k'.

→ It is an Non-Negative Dist<sup>n</sup>.

→ It is an Right Skewed Dist<sup>n</sup>.



⇒ Chi-Square Test ( $\chi^2$ ) =  $\sum (O - E)^2 / E$

→ The Chi-Square Test is Based on Chi-Square dist<sup>n</sup>.

→ Goodness of fit Test is used to perform this.

$$\chi^2_{\text{stat.}} = \sum \frac{(O - E)^2}{E}$$

Expected

Q → 12% of people are left Handed. To Verify The Theory, We Took an Sample of 75 students where 11 were left Handed. ( $\alpha = 0.05$ )

→ Expected →  $\text{Left - Handed} = 12\% \text{ of } 75 = 9$   
 $\text{Right - Handed} = 75 - 9 = 66$

Observed →  $\text{Left - Handed} = 11$   
 $\text{Right Handed} = 75 - 11 = 64$

	Expected (E)	Observed (O)
Left	9	11
Right	66	64

① Examine The Hypothesis →  $H_0: \mu = 12\%$  vs  $H_A: \mu \neq 12\%$

② Level of Sig ( $\alpha$ ) = 0.05 (One Tailed Test)  
~~(No 2-Tailed Test Possible In  $\chi^2$ -Test)~~

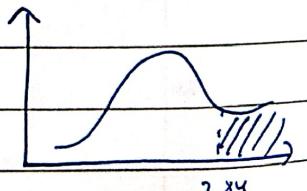
③ Type of Test =  $\chi^2$ -Test

④  $\chi^2_{\text{std.}} = \sum \frac{(O-E)^2}{E} = \frac{(11-9)^2}{9} + \frac{(64-66)^2}{66} = 0.505$

Now,  $\alpha = 0.05$ , dof = 1  $\Rightarrow \chi^2_{\text{critic}} = 3.84$

$\chi^2_{\text{critic}} (\alpha = 0.05, \text{dof} = 1) = 3.84$

∴ Now, As,  $\chi^2_{\text{std.}} = 0.505 < 3.84$



∴ Thus, We fail To Reject The Null Hypothesis.

12% of People are left Handed with 95% Confidence.



Q → For a sample of 500 people, observed & Expected values were :-

	Expected (E)	observed (o)
$W \leq 50 \text{ kg}$	$20\% = 100$	140
$50 < W < 75$	$30\% = 150$	160
$W \geq 75 \text{ kg}$	$50\% = 250$	200

→ Using  $\alpha = 0.05$ , check the expected values are correct?

① → Forme The Hypothesis  $\rightarrow$   $H_0$  : The data is as per expectation  
 $\rightarrow$   $H_A$  : " " " Not " "

② Level of Sig ( $\alpha$ ) = 0.05

③ Type of Test :  $\chi^2$  - Test

$$\begin{aligned} \text{④ } \chi^2_{\text{stat.}} &= \sum \frac{(O-E)^2}{E} = \frac{(140-100)^2}{100} + \frac{(160-150)^2}{150} + \frac{(200-250)^2}{250} \\ &= \underline{\underline{26.66}} \end{aligned}$$

⑤ Now,  $\alpha = 0.05$  and  $\text{dof} = n-1 = 3-1 = 2$

Thus,  $\chi^2_{\text{critise}} (\alpha=0.05, \text{dof}=2) = 5.99$

∴ Now,  $\chi^2_{\text{stat.}} = 26.66 \gg \chi^2_{\text{critise}}$

∴ Thus, Reject The Null Hypothesis.

∴ The observed values of weight are ~~not~~ changed as were expected.

### Chi-square Distribution Table

d.f.	.995	.99	.975	.95	.9	.1	.05	.025	.01
1	0.00	0.00	0.00	0.00	0.02	2.71	3.84	5.02	6.63
2	0.01	0.02	0.05	0.10	0.21	4.61	5.99	7.38	9.21
3	0.07	0.11	0.22	0.35	0.58	6.25	7.81	9.35	11.34
4	0.21	0.30	0.48	0.71	1.06	7.78	9.49	11.14	13.28
5	0.41	0.55	0.83	1.15	1.61	9.24	11.07	12.83	15.09
6	0.68	0.87	1.24	1.64	2.20	10.64	12.59	14.45	16.81
7	0.99	1.24	1.69	2.17	2.83	12.02	14.07	16.01	18.48
8	1.34	1.65	2.18	2.73	3.49	13.36	15.51	17.53	20.09
9	1.73	2.09	2.70	3.33	4.17	14.68	16.92	19.02	21.67
10	2.16	2.56	3.25	3.94	4.87	15.99	18.31	20.48	23.21
11	2.60	3.05	3.82	4.57	5.58	17.28	19.68	21.92	24.72
12	3.07	3.57	4.40	5.23	6.30	18.55	21.03	23.34	26.22
13	3.57	4.11	5.01	5.89	7.04	19.81	22.36	24.74	27.69
14	4.07	4.66	5.63	6.57	7.79	21.06	23.68	26.12	29.14
15	4.60	5.23	6.26	7.26	8.55	22.31	25.00	27.49	30.58
16	5.14	5.81	6.91	7.96	9.31	23.54	26.30	28.85	32.00
17	5.70	6.41	7.56	8.67	10.09	24.77	27.59	30.19	33.41
18	6.26	7.01	8.23	9.39	10.86	25.99	28.87	31.53	34.81
19	6.84	7.63	8.91	10.12	11.65	27.20	30.14	32.85	36.19
20	7.43	8.26	9.59	10.85	12.44	28.41	31.41	34.17	37.57
22	8.64	9.54	10.98	12.34	14.04	30.81	33.92	36.78	40.29
24	9.89	10.86	12.40	13.85	15.66	33.20	36.42	39.36	42.98
26	11.16	12.20	13.84	15.38	17.29	35.56	38.89	41.92	45.64
28	12.46	13.56	15.31	16.93	18.94	37.92	41.34	44.46	48.28
30	13.79	14.95	16.79	18.49	20.60	40.26	43.77	46.98	50.89
32	15.13	16.36	18.29	20.07	22.27	42.58	46.19	49.48	53.49
34	16.50	17.79	19.81	21.66	23.95	44.90	48.60	51.97	56.06
38	19.29	20.69	22.88	24.88	27.34	49.51	53.38	56.90	61.16
42	22.14	23.65	26.00	28.14	30.77	54.09	58.12	61.78	66.21
46	25.04	26.66	29.16	31.44	34.22	58.64	62.83	66.62	71.20
50	27.99	29.71	32.36	34.76	37.69	63.17	67.50	71.42	76.15
55	31.73	33.57	36.40	38.96	42.06	68.80	73.31	77.38	82.29
60	35.53	37.48	40.48	43.19	46.46	74.40	79.08	83.30	88.38
65	39.38	41.44	44.60	47.45	50.88	79.97	84.82	89.18	94.42
70	43.28	45.44	48.76	51.74	55.33	85.53	90.53	95.02	100.43
75	47.21	49.48	52.94	56.05	59.79	91.06	96.22	100.84	106.39
80	51.17	53.54	57.15	60.39	64.28	96.58	101.88	106.63	112.33
85	55.17	57.63	61.39	64.75	68.78	102.08	107.52	112.39	118.24
90	59.20	61.75	65.65	69.13	73.29	107.57	113.15	118.14	124.12
95	63.25	65.90	69.92	73.52	77.82	113.04	118.75	123.86	129.97
100	67.33	70.06	74.22	77.93	82.36	118.50	124.34	129.56	135.81

## ⇒ Baye's Theorem

→ Probability is the chance of success / # possible outcomes.

# The Rules of Probability are :-

- For any Event,  $0 \leq P(A) \leq 1$
- $\sum \text{Prob.} = 1$ .
- $P(\bar{A}) = 1 - P(A)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- The Multiplication Rule for Prob. Is :-

$$P(A \cap B) = P(A) \times P(B/A) \quad \text{--- (1)}$$

$$P(B \cap A) = P(B) \times P(A/B) \quad \text{--- (2)}$$

Then,  $P(A) \times P(B/A) = P(B) \times P(A/B)$

$$P(A/B) = \frac{P(A) \times P(B/A)}{P(B)}$$

⇒ 10% of Patients in an clinic have Liver disease. 25% of the Clinic patients are Alcoholic. Among the Healed Patients, 7% are Alcoholic.

→ What is the prob. of patients having liver disease given that he is Alcoholic?

→ Let, A = liver disease & B = patient is Alcoholic

$$P(A) = 0.10, \quad P(B) = 0.05, \quad P(B/A) = 0.07$$

$$\therefore P(A/B) = \frac{P(A) \times P(B/A)}{P(B)} = \frac{0.10 \times 0.07}{0.05} = 0.14 = 14\%$$

∴ Thus, 14% chance that Patient detected with Liver disease was Alcoholic.



→ F - Distribution (Fisher - Snedecor distn)

→ The F - Dist<sup>n</sup> is an Probability dist<sup>n</sup> that has been useful in context of Comparing Variances of Two or more Samples.

→ It is Right Skewed and Takes Only Non-Negative Values.

→ The F - dist<sup>n</sup> with d<sub>1</sub> and d<sub>2</sub> degree of freedom is the distn :-

$$X = \frac{S_1^2 / \text{dof}_1}{S_2^2 / \text{dof}_2}$$

$$\frac{S_1^2 / \text{dof}_1}{S_2^2 / \text{dof}_2} \sim F_{\text{dof}_1, \text{dof}_2}$$

P

f(x)

CH

←

12

10

→ F statistic  $F = \frac{S_1^2}{S_2^2}$  (Variance Ratio Test) (No Two Tailed Test)

② Where, S<sub>1</sub> and S<sub>2</sub> are Sample Std Deviation.

→ The Following data is about the # Bulbs produced daily by 2 workers :-

$$A \quad B \quad (\alpha = 0.05)$$

$$40 \quad 39 \quad \frac{(\bar{x} - x)}{\sqrt{s^2}} = -0.2$$

$$30 \quad 38 \quad 1.0$$

$$38 \quad 41$$

$$41 \quad \bar{x} \quad 33.5 \quad \frac{(\bar{x} - x)}{\sqrt{s^2}} = 2.47$$

$$38 \quad 32.5$$

$$35 \quad 39.5$$

$$31 \quad 40.5$$

$$31 \quad 34.5$$

$$28 \quad 35.5$$

② Can we conclude based on 'B' that Worker B is more stable & effective?

$$P \quad f(x) \quad CH$$

★ Mean can't be used for this Test because mean is almost same for both samples. Thus, we will compare variance.



① Frame The Hypothesis  $\rightarrow H_0 : S_1^2 = S_2^2$

(null)  $\rightarrow H_A : S_1^2 \neq S_2^2$

② Level of Sig ( $\alpha$ ) = 0.05

After year for test to be carried out

③ Type of Test:  $\rightarrow$  F-test for homogeneity of variance

With  $n_1 = 8$  &  $n_2 = 7$  degrees of freedom

④ Firstly Compute  $S_1^2$  and  $S_2^2$ ,  $\bar{x}_1 = 37$

$$x_1 - \bar{x}_1 \quad (x_1 - \bar{x}_1)^2$$

$$\text{For } S_1 \rightarrow 40 \quad 37 \quad 9$$

$$(30.0 - 37) \quad (30.0 - 37) \quad 49$$

$$38 \quad 37 \quad 1$$

$$\therefore S_1^2 = \frac{\sum (x_1 - \bar{x}_1)^2}{n-1} = \frac{80}{7-1} = 16$$

$$\text{For } S_2 \rightarrow \bar{x}_2 = 37 \quad 14 \quad (x_2 - \bar{x}_2)^2$$

$$39 \quad 37 \quad 82 \quad 4$$

$$38 \quad 37 \quad 28 \quad 1$$

$$41 \quad 37 \quad 16$$

$$33 \quad 37 \quad 16$$

$$32 \quad 37 \quad 25$$

$$40 \quad 37 \quad 9$$

$$34 \quad 37 \quad 9$$

$$\therefore S_2^2 = \frac{\sum (x_2 - \bar{x}_2)^2}{n-1} = \frac{84}{8-1} = 12$$



$$\therefore F_{\text{stat.}} = \frac{s_1^2}{s_2^2} = \frac{16}{12} = \underline{\underline{1.33}}$$

(5) F<sub>critical</sub>,  $\alpha = 0.05$ ,  $\text{dof}_1 = n_A - 1 = 6 - 1 = 5$   
 $\text{dof}_2 = n_B - 1 = 8 - 1 = 7$

$\therefore$  Thus, By F - Table,

$$\text{F}_{\text{critical}} (\alpha = 0.05, \text{dof} = (5, 7)) = \underline{\underline{3.97}}$$

$\therefore$  Thus, we Fail To Reject The Null Hypothesis.

$\therefore$  Worker B Is Not more Effective Than Worker A.

## F distribution critical value landmarks

Table entries are critical values for  $F^*$  with probably  $p$  in right tail of the distribution.

Figure of F distribution (like in Moore, 2004, p. 656) here.

Degrees of freedom in denominator (df2)	p	Degrees of freedom in numerator (df1)										
		1	2	3	4	5	6	7	8	12	24	1000
1	0.100	39.86	49.50	53.59	55.83	57.24	58.20	58.91	59.44	60.71	62.00	63.30
	0.050	161.4	199.5	215.7	224.6	230.2	234.0	236.8	238.9	243.9	249.1	254.2
	0.025	647.8	799.5	864.2	899.6	921.8	937.1	948.2	956.6	976.7	997.3	1017.8
	0.010	4052	4999	5404	5624	5764	5859	5928	5981	6107	6234	6363
	0.001	405312	499725	540257	562668	576496	586033	593185	597954	610352	623703	636101
2	0.100	8.53	9.00	9.16	9.24	9.29	9.33	9.35	9.37	9.41	9.45	9.49
	0.050	18.51	19.00	19.16	19.25	19.30	19.33	19.35	19.37	19.41	19.45	19.49
	0.025	38.51	39.00	39.17	39.25	39.30	39.33	39.36	39.37	39.41	39.46	39.50
	0.010	98.50	99.00	99.16	99.25	99.30	99.33	99.36	99.38	99.42	99.46	99.50
	0.001	998.38	998.84	999.31	999.31	999.31	999.31	999.31	999.31	999.31	999.31	999.31
3	0.100	5.54	5.46	5.39	5.34	5.31	5.28	5.27	5.25	5.22	5.18	5.13
	0.050	10.13	9.55	9.28	9.12	9.01	8.94	8.89	8.85	8.74	8.64	8.53
	0.025	17.44	16.04	15.44	15.10	14.88	14.73	14.62	14.54	14.34	14.12	13.91
	0.010	34.12	30.82	29.46	28.71	28.24	27.91	27.67	27.49	27.05	26.60	26.14
	0.001	167.06	148.49	141.10	137.08	134.58	132.83	131.61	130.62	128.32	125.93	123.52
4	0.100	4.54	4.32	4.19	4.11	4.05	4.01	3.98	3.95	3.90	3.83	3.76
	0.050	7.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04	5.91	5.77	5.63
	0.025	12.22	10.65	9.98	9.60	9.36	9.20	9.07	8.98	8.75	8.51	8.26
	0.010	21.20	18.00	16.69	15.98	15.52	15.21	14.98	14.80	14.37	13.93	13.47
	0.001	74.13	61.25	56.17	53.43	51.72	50.52	49.65	49.00	47.41	45.77	44.09
5	0.100	4.06	3.78	3.62	3.52	3.45	3.40	3.37	3.34	3.27	3.19	3.11
	0.050	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82	4.68	4.53	4.37
	0.025	10.01	8.43	7.76	7.39	7.15	6.98	6.85	6.76	6.52	6.28	6.02
	0.010	16.26	13.27	12.06	11.39	10.97	10.67	10.46	10.29	9.89	9.47	9.03
	0.001	47.18	37.12	33.20	31.08	29.75	28.83	28.17	27.65	26.42	25.13	23.82
6	0.100	3.78	3.46	3.29	3.18	3.11	3.05	3.01	2.98	2.90	2.82	2.72
	0.050	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15	4.00	3.84	3.67
	0.025	8.81	7.26	6.60	6.23	5.99	5.82	5.70	5.60	5.37	5.12	4.86
	0.010	13.75	10.92	9.78	9.15	8.75	8.47	8.26	8.10	7.72	7.31	6.89
	0.001	35.51	27.00	23.71	21.92	20.80	20.03	19.46	19.03	17.99	16.90	15.77
7	0.100	3.59	3.26	3.07	2.96	2.88	2.83	2.78	2.75	2.67	2.58	2.47
	0.050	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73	3.57	3.41	3.23
	0.025	8.07	6.54	5.89	5.52	5.29	5.12	4.99	4.90	4.67	4.41	4.15
	0.010	12.25	9.55	8.45	7.85	7.46	7.19	6.99	6.84	6.47	6.07	5.66
	0.001	29.25	21.69	18.77	17.20	16.21	15.52	15.02	14.63	13.71	12.73	11.72
8	0.100	3.46	3.11	2.92	2.81	2.73	2.67	2.62	2.59	2.50	2.40	2.30
	0.050	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44	3.28	3.12	2.93
	0.025	7.57	6.06	5.42	5.05	4.82	4.65	4.53	4.43	4.20	3.95	3.68
	0.010	11.26	8.65	7.59	7.01	6.63	6.37	6.18	6.03	5.67	5.28	4.87
	0.001	25.41	18.49	15.83	14.39	13.48	12.86	12.40	12.05	11.19	10.30	9.36
9	0.100	3.36	3.01	2.81	2.69	2.61	2.55	2.51	2.47	2.38	2.28	2.16
	0.050	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23	3.07	2.90	2.71
	0.025	7.21	5.71	5.08	4.72	4.48	4.32	4.20	4.10	3.87	3.61	3.34
	0.010	10.56	8.02	6.99	6.42	6.06	5.80	5.61	5.47	5.11	4.73	4.32
	0.001	22.86	16.39	13.90	12.56	11.71	11.13	10.70	10.37	9.57	8.72	7.84

Critical values computed with Excel 9.0

		Degrees of freedom in numerator (df1)										
		1	2	3	4	5	6	7	8	12	24	1000
Degrees of freedom in denominator (df2)	10	0.100	3.29	2.92	2.73	2.61	2.52	2.46	2.41	2.38	2.28	2.18
	10	0.050	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07	2.91	2.74
	10	0.025	6.94	5.46	4.83	4.47	4.24	4.07	3.95	3.85	3.62	3.37
	10	0.010	10.04	7.56	6.55	5.99	5.64	5.39	5.20	5.06	4.71	4.33
	10	0.001	21.04	14.90	12.55	11.28	10.48	9.93	9.52	9.20	8.45	7.64
	12	0.100	3.18	2.81	2.61	2.48	2.39	2.33	2.28	2.24	2.15	2.04
	12	0.050	4.75	3.89	3.49	3.26	3.11	3.00	2.91	2.85	2.69	2.51
	12	0.025	6.55	5.10	4.47	4.12	3.89	3.73	3.61	3.51	3.28	3.02
	12	0.010	9.33	6.93	5.95	5.41	5.06	4.82	4.64	4.50	4.16	3.78
	12	0.001	18.64	12.97	10.80	9.63	8.89	8.38	8.00	7.71	7.00	6.25
Degrees of freedom in denominator (df2)	14	0.100	3.10	2.73	2.52	2.39	2.31	2.24	2.19	2.15	2.05	1.94
	14	0.050	4.60	3.74	3.34	3.11	2.96	2.85	2.76	2.70	2.53	2.35
	14	0.025	6.30	4.86	4.24	3.89	3.66	3.50	3.38	3.29	3.05	2.79
	14	0.010	8.86	6.51	5.56	5.04	4.69	4.46	4.28	4.14	3.80	3.43
	14	0.001	17.14	11.78	9.73	8.62	7.92	7.44	7.08	6.80	6.13	5.41
	16	0.100	3.05	2.67	2.46	2.33	2.24	2.18	2.13	2.09	1.99	1.87
	16	0.050	4.49	3.63	3.24	3.01	2.85	2.74	2.66	2.59	2.42	2.24
	16	0.025	6.12	4.69	4.08	3.73	3.50	3.34	3.22	3.12	2.89	2.63
	16	0.010	8.53	6.23	5.29	4.77	4.44	4.20	4.03	3.89	3.55	3.18
	16	0.001	16.12	10.97	9.01	7.94	7.27	6.80	6.46	6.20	5.55	4.85
Degrees of freedom in denominator (df2)	18	0.100	3.01	2.62	2.42	2.29	2.20	2.13	2.08	2.04	1.93	1.81
	18	0.050	4.41	3.55	3.16	2.93	2.77	2.66	2.58	2.51	2.34	2.15
	18	0.025	5.98	4.56	3.95	3.61	3.38	3.22	3.10	3.01	2.77	2.50
	18	0.010	8.29	6.01	5.09	4.58	4.25	4.01	3.84	3.71	3.37	3.00
	18	0.001	15.38	10.39	8.49	7.46	6.81	6.35	6.02	5.76	5.13	4.45
	20	0.100	2.97	2.59	2.38	2.25	2.16	2.09	2.04	2.00	1.89	1.77
	20	0.050	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45	2.28	2.08
	20	0.025	5.87	4.46	3.86	3.51	3.29	3.13	3.01	2.91	2.68	2.41
	20	0.010	8.10	5.85	4.94	4.43	4.10	3.87	3.70	3.56	3.23	2.86
	20	0.001	14.82	9.95	8.10	7.10	6.46	6.02	5.69	5.44	4.82	4.15
Degrees of freedom in denominator (df2)	30	0.100	2.88	2.49	2.28	2.14	2.05	1.98	1.93	1.88	1.77	1.64
	30	0.050	4.17	3.32	2.92	2.69	2.53	2.42	2.33	2.27	2.09	1.89
	30	0.025	5.57	4.18	3.59	3.25	3.03	2.87	2.75	2.65	2.41	2.14
	30	0.010	7.56	5.39	4.51	4.02	3.70	3.47	3.30	3.17	2.84	2.47
	30	0.001	13.29	8.77	7.05	6.12	5.53	5.12	4.82	4.58	4.00	3.36
	50	0.100	2.81	2.41	2.20	2.06	1.97	1.90	1.84	1.80	1.68	1.54
	50	0.050	4.03	3.18	2.79	2.56	2.40	2.29	2.20	2.13	1.95	1.74
	50	0.025	5.34	3.97	3.39	3.05	2.83	2.67	2.55	2.46	2.22	1.93
	50	0.010	7.17	5.06	4.20	3.72	3.41	3.19	3.02	2.89	2.56	2.18
	50	0.001	12.22	7.96	6.34	5.46	4.90	4.51	4.22	4.00	3.44	2.82
Degrees of freedom in denominator (df2)	100	0.100	2.76	2.36	2.14	2.00	1.91	1.83	1.78	1.73	1.61	1.46
	100	0.050	3.94	3.09	2.70	2.46	2.31	2.19	2.10	2.03	1.85	1.63
	100	0.025	5.18	3.83	3.25	2.92	2.70	2.54	2.42	2.32	2.08	1.78
	100	0.010	6.90	4.82	3.98	3.51	3.21	2.99	2.82	2.69	2.37	1.98
	100	0.001	11.50	7.41	5.86	5.02	4.48	4.11	3.83	3.61	3.07	2.46
	1000	0.100	2.71	2.31	2.09	1.95	1.85	1.78	1.72	1.68	1.55	1.39
	1000	0.050	3.85	3.00	2.61	2.38	2.22	2.11	2.02	1.95	1.76	1.53
	1000	0.025	5.04	3.70	3.13	2.80	2.58	2.42	2.30	2.20	1.96	1.65
	1000	0.010	6.66	4.63	3.80	3.34	3.04	2.82	2.66	2.53	2.20	1.81
	1000	0.001	10.89	6.96	5.46	4.65	4.14	3.78	3.51	3.30	2.77	2.16

Use StaTable, WinPepi > WhatIs, or other reliable software to determine specific  $p$  values

## ⇒ ANOVA (Analysis of Variance) & It's Assumptions

- ANOVA is a Statistical method used To Compare The means of 2 + Groups.
- Eg:- Effect of 0mg, 10mg & 100mg of medicine.

### # Assumptions of ANOVA

- The Population from which Samples are drawn should be normally distributed.
- The Sample should be Independent from each other.
- There should be The Absence of The outliers.
- Homogeneity of Variance : Homogeneity means That The Variance among The Groups should be Approximately Equal.

$$\text{Eg } S_1, S_2, S_3 \Rightarrow S_1^2 = S_2^2 = S_3^2$$



## ⇒ Types of ANOVA

1. One Way ANOVA → One Factor with Atleast Two levels and The levels are been Independent.

Eg:  $F = \frac{SST}{SSE} = \frac{1 - M}{M}$

Medicine

0 mg

10 mg

100 mg (dosage)

2. Repeated Measures ANOVA → One Factor with Atleast Two levels but The levels are different.

Eg: - A person day 1 works 2 hrs, day 2 works 3 hrs & day 4 works 4 hrs  
GYM 2 hrs 1 hr 1/2 hr 1/2 hr  
(my day workout depend on free. day)

3. Factorial ANOVA → Two or more factors (with 2 + levels) - levels

can be either dependent, Independent or Both.

→ It Is Also Known as 2 Way ANOVA as IT deals with 2 factors.

## ⇒ ANOVA Test

→ Hypothesis Testing In ANOVA (Partitioning of Variance In ANOVA)

$H_0: \mu_1 = \mu_2 = \dots = \mu_k$

$H_A: \text{Atleast one of The Sample mean is Not Equal}$

→ Test Statistics =  $\frac{\text{Variance b/w Samples}}{\text{Variance within Samples}}$

Variance within Samples

Q → There are 3 dosage of an medicine given to three samples of the patients. Rating is done if medicine works (1 - 10). Are there any differences in the 3 cond'n ( $\alpha = 0.05$ )

	<u>(S<sub>1</sub>)</u>	<u>(S<sub>2</sub>)</u>	<u>(S<sub>3</sub>)</u>
0 mg	long	100 mg	
+ 4	7	4	
8 (i)	6 P.S.O. = mean 3		
7	6	2	
8	7	3	
8	8 with one more 4 and (ii)		
9	7.2	3	
8	6	2	
<u>57</u>	<u>47</u>	<u>21</u>	

① → Frame The Hypothesis :  $H_0 \rightarrow \mu_{S_1} = \mu_{S_2} = \mu_{S_3}$

$H_A$  → Any of The mean Is Not Equal

(2) Level of Sig. ( $\alpha$ ) = 0.05 (One Tailed Test)  
(No Two-Tailed present In ANOVA-Test)

## Type of Test : ANOVA - Test

(4) ~~F statistics~~ =  $\frac{\text{Variance b/w Samples}}{\text{Variance b/w Samples}}$

Variance Within Samples defined as  $s^2_w$

Rs. 0.00 + 81 = Rs. 0.01. *Original price*

Sum of Squares (SS)      df      SS/df

B/W Samples 2008 - 2011 = 2011

Within Samples ~~mm~~ 2.3 ~~mm~~ 2.4 ~~mm~~

Total



(i) Sum of Squares Within Groups / Samples

$$SS_{\text{within}} = \sum y_i^2 - n \frac{(\sum a_i)^2}{n}$$

$$SS_{\text{within}} = (9^2 + 8^2 + 7^2 + \dots + 7^2 + 6^2 + 6^2 + \dots + 4^2 + 3^2 + 2^2 + \dots + 3^2 + 2^2)$$

$$= \left( \frac{57^2 + 47^2 + 21^2}{7} \right)$$

$$SS_{\text{within}} = 10.29 \quad \text{--- (1)}$$

(ii) Sum of Squares B/W Samples

$$SS_{B/W} = \frac{\sum (\sum a_i)^2}{n} - \frac{T^2}{N} \quad (T \equiv \text{Total Sum})$$

$$SS_{B/W} = (57^2 + 47^2 + 21^2) - \frac{(57 + 47 + 21)^2}{21}$$

$$SS_{B/W} = 48.67 \quad \text{--- (2)}$$

$$(iii) \text{ dof}_{B/W \text{ Samples}} = (\# \text{ Samples}) - 1 = 3 - 1 = 2 \quad \text{--- (3)}$$

$$(iv) \text{ dof}_{\text{within samples}} = (\# \text{ Total Elements}) - (\# \text{ Samples}) = 21 - 3 = 18 \quad \text{--- (4)}$$

$\therefore$  Sum of Sq (SS) =  $\text{dof} \times \text{M.S.} (SS/\text{dof})$

$$\text{B/W The Samples} \quad 48.67 \quad 2 \quad 49.34$$

$$\text{Within Samples} \quad 10.29 \quad 18 \quad 0.54$$

$$108.96 \quad 20$$

$$F_{\text{stat}} = \frac{\text{MS}_{\text{Between}}}{\text{MS}_{\text{within}}} = \frac{49.34}{0.54} = 86.56$$



(5)

$$F_{\text{critical}} \text{ for } \alpha = 0.05, \quad \text{dof}_{B/W} = (\#\text{Samples}) - 1 = 3 - 1 = 2 \\ \text{dof}_{\text{within}} = N - (\#\text{Samples}) = 21 - 3 = 18$$

∴ By F-Table,  $F_{\text{critical}} (\alpha = 0.05, \text{dof} = (2, 18)) = 3.55$

∴ As,  $F_{\text{stat.}} = 86.56 >> 3.55$

∴ Thus, Reject The Null Hypothesis.

∴ Thus, Any of The mean Is Not Equal.

— JAI LOGIC